

## Homework 4 – Due 5 October

The total points on this homework is 175. Out of these, 5 points are reserved for clarity of presentation, punctuation and commenting with respect to the code.

The homework for this week are mainly exercises in using the `apply()` function, and its benefits in many cases where the problem can be “reduced” (which may mean, expanded for some cases) to operations at the margins of an array. Some of the problems are for a bigger dataset that is more cumbersome to notice operations on, so you are advised to try it out first on a small array (say of dimension  $3 \times 4 \times 5$ ) and then moving to the given problem once you are confident of your approach to solving the problem.

1. The file `fbp-img.dat` on Canvas is a  $128 \times 128$  matrix containing the results of an image of fluorodeoxyglucose (FDG-18) radio-tracer intake reconstructed using Positron Emission Tomography (PET).
  - (a) Read in the data as a matrix. [5 points]
  - (b) Use the `image()` function in R with your choice of color map to image the data. [5 points]  
*(At this point, it may be worth understanding how the `image()` function does the plotting – it plots the columns of the matrix from left to right (which may be what we would like) but the rows of the matrix from bottom to top (which may not be what we would like). This issue did not arise in the previous question because I had changed the columns to be bottoms-up so that your plot turned out to be the way we would prefer it. In this example, you should use index vectors to reverse the display so that you get the correct-looking version of the brain. Recall that in most situations, eg. scan that you would get at a clinic, the front of a person’s brain/head is imaged facing the top while the back of the head is towards the bottom.)*
  - (c) We will now (albeit, somewhat crudely) compress the data in the following two ways:
    - i. Obtain the range of values in the matrix, subdividing into 16 bins. Bin the values in the matrix and replace each value with the mid-point of each bin. Image these new binned matrix values. [10 points]
    - ii. In this second case, we will group according to the 16 equally-spaced quantile bins, which can be obtained using the `quantile()` function in R with appropriate arguments. Use these obtained quantile bins to group the data, replacing each entry in the matrix with its mid-point. Image these new binned matrix values. [8 points]
    - iii. Comment, and discuss similarities and differences on the differences in the two images thus obtained with the original image. [2 points]
2. *Microarray gene expression data.* The file, accessible on Canvas at `diurnaldata.csv` contains gene expression data on 22,810 genes from Arabidopsis plants exposed to equal periods of light and darkness in the diurnal cycle. Leaves were harvested at eleven time-points, at the start of the experiment (end of the light period) and subsequently after 1, 2, 4, 8 and 12 hours of darkness and light each. Note that there are 23 columns, with the first column representing the gene probeset. Columns 2–12 represent measurements on gene abundance taken at 1, 2, 4, 8 and 12 hours of darkness and light each, while columns 13–23 represent the same for a second replication.
  - (a) Read in the dataset. Note that this is a big file, and can take a while, especially on a slow connection. [5 points]
  - (b) For each gene, calculate the mean abundance level at each time-point, and store the result in a matrix. One way to achieve this is to create a three-dimensional array of dimension  $22,810 \times 11 \times 2$  and to use `apply` over it with the appropriate function and over the appropriate margins. Note that you are only asked present the commands that you use here. From now on, we will use this dataset averaged over the two replications for each gene. [15 points]
  - (c) *Standardization.* Gene data are compared to each other by way of correlations. (Correlation between any two sequences is related, by means of an affine (linear) transformation, to the Euclidean distance between the sequences, after standardization to have mean zero and standard deviation 1).

- i. In pursuance of the objective, use the `apply` function to calculate the mean of the mean abundance level over all time-points for each gene. [10 points]
    - ii. Set up a matrix of dimension  $22,810 \times 11$  of the means of the mean abundance levels, where each row is a replicated version of the first one. Use this to eliminate the mean effect from the matrix stored in the previous part. [10 points]
    - iii. Use the `apply` function again to calculate the standard deviation of each row of the matrix in the part (b) above, and proceed to obtain the scaled measurements on the genes. [10 points]
  - (d) The file `micromeans.dat` contains a  $20 \times 11$  matrix of measurements. Read in this dataset, and standardize as above. [5 points]
  - (e) We will now identify which of these means is closest to each gene. To do so, set up a three-dimensional array (of dimension  $22,810 \times 11 \times 20$ ) with 20 replicated datasets (of the 22,810 diurnal mean abundance levels over 11 time-points). Next, set up replications of the 20 mean vectors into a three-dimensional array (of dimension  $22,810 \times 11 \times 20$ ), with the replications occurring along the first dimension. Using `apply` and the above, obtain a  $22,810 \times 20$  matrix of Euclidean distances for each gene and mean. Finally, obtain the means to which each gene is closest to. Report the frequency distribution, and tabulate the frequencies, and display using a piechart. [25 points]
3. The United States Postal Service has had a long-term project to automating the recognition of handwritten digits for zip codes. The file `ziptrain.dat` has data on different numbers of specimens for each digit. Each observation is in the form of a 256-dimensional vector of pixel intensities. These form a  $16 \times 16$  image of pixel intensities for each digit. Each digit is determined to be a specimen of the actual digit given in the file `zipdigit.dat`. The objective is to clean up the dataset to further distinguish one digit from the other.
- (a) Read in the dataset as a vector. This dataset has 2000 rows and 256 columns. Convert and store the dataset as a three-dimensional array of dimensions  $16 \times 16 \times 2000$ . [5 points]
  - (b) Our objective is to image all the 2000 pictures in a  $40 \times 50$  display. We will hasten this display process by setting up the stage for using `apply`.
    - i. If we image the second record, we note that without any changes, the image mirrors the digit “5”. Fix these by reversing the appropriate dimension of the array. Also note that the white space around the axes is wasted space and should be removed. To substantially reduce the white space, consider `par(mar = rep(0.05, 4))`. Further, the axes conveys no information so consider its elimination, using `axes = F` in the call to the `image()` function. Finally, we may as well use a gray scale here so consider using `col = gray(31:0/31)` in the image argument. Display the second record as an image, after addressing all these issues. [5 points]
    - ii. Using `par(mfrow = c(40, 50))` and a total image size of  $5.2'' \times 6.5''$ , image all the 2000 digits using `apply` and the experience you gained in part i. above. [10 points]
  - (c) We now compute the mean and standard deviation images of the digits. To do so, we convert the array back to a matrix, and calculate the means for the ten digits, convert back to a three-dimensional array. Do exactly that. Also, compute and display the standard deviation images in exactly a similar way. [20 points]
  - (d) Our final act will involve cleaning up the dataset. To do so, remove the mean digit effect from each record. Then, perform a singular value decomposition on the  $2000 \times 256$  matrix of deviations  $\mathbf{Y}$  (that is the records with the corresponding digit mean effect removed). (Recall that the SVD is  $\mathbf{Y} = \mathbf{U}\mathbf{D}\mathbf{V}'$  and is obtained via the `svd()` function in R.) Replace  $\mathbf{D}$  with a diagonal matrix  $\mathbf{D}_k$  of the first  $k$  eigenvalues of  $\mathbf{D}$  and the remainder as zero. Then let  $\mathbf{Y}_k = \mathbf{U}\mathbf{D}_k\mathbf{V}'$ . Add back the mean and display the resulting images of all the digits for  $k = 25, 50, 75$  in the same manner as (b)ii. Comment on the displays.<sup>1</sup> [30 points]

<sup>1</sup>The reduced representation image of the original values is similar to keeping only the lower frequencies which keeps the main definitions and features of the images but filters out the noise. This is studied further in Stat 5010, among other classes. You will also see a version of this in the department seminar on October 7, 2024.