

STAT 3005 – Nonparametric Statistics

ASSIGNMENT 4

Due: 15 November (Friday) at 1800

Fall 2024

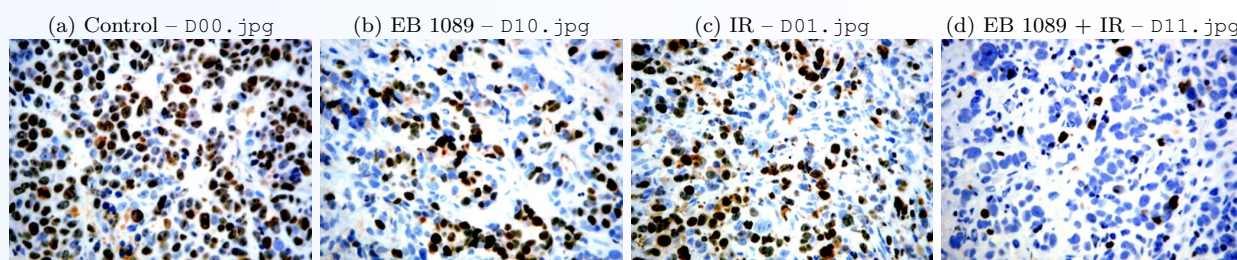
Exercise 4.1 (★★☆ — Cancer image data analysis (80%)). In RGB format, a color image $C = (C_{ijk}) \in \mathbb{R}^{n_1 \times n_2 \times 3}$ can be represented by a tensor, i.e., a multidimensional array, of $n_1 \times n_2$ pixel values of the three primary colors (red, green and blue). Each value $C_{i,j,k}$ is rescaled to $[0, 1]$ so that it represents the intensity of the k th primary color at the (i, j) pixel. For example, the (i, j) th pixel is the triple $(C_{i,j,1}, C_{i,j,2}, C_{i,j,3})$.

- It is black if $C_{i,j,1:3} = (0, 0, 0)$.
- It is white if $C_{i,j,1:3} = (1, 1, 1)$.
- It is blue if $C_{i,j,1:3} = (0, 0, 1)$.
- It is yellow if $C_{i,j,1:3} = (1, 1, 0)$.
- It is purple if $C_{i,j,1:3} = (0.5, 0, 0.5)$.

A greyscale (monochrome) image can be formed by a linear combination of the intensities of the three primary colors. Typically, the greyscale version of C is a matrix $G = (G_{i,j}) \in \mathbb{R}^{n_1 \times n_2}$ defined as

$$G_{i,j} = 0.22C_{i,j,1} + 0.71C_{i,j,2} + 0.07C_{i,j,3}, \quad (4.1)$$

for $i = 1, \dots, n_1$ and $j = 1, \dots, n_2$. The closer $G_{i,j}$ is to zero, the darker is the (i, j) th pixel. In this exercise, we are interested in performing histological analysis of cancer treatment. In particular, we compare images of tumor tissue section published in Clinical Cancer Research by Sundaram et al. (2003); see this article ([click here](#)). The following images show the proliferative activity of tumors exposed to (a) control, (b) EB 1089, (c) radiation, and (d) EB 1089 + radiation by Ki67 staining of paraffin-embedded tumor tissue sections. The dark regions indicate cancer cells. The images can be downloaded [here](#). Our goal is to test if there is any cancer treatment effect.



- (20%) Complete the following R function `get.gray` for producing grayscale images according to (4.1). Use this function to produce grayscale version of the above 4 images. Plot them.

```

1 #install.packages("jpeg")
2 img = rep(list(NA), 4)
3 img[[1]] = jpeg::readJPEG("D00.jpg")
4 img[[2]] = jpeg::readJPEG("D10.jpg")
5 img[[3]] = ... # complete this line
6 img[[4]] = ... # complete this line
7 # plots of color images
8 par(mfrow=c(1,4),mar=c(0,0,1,0))
9 for(i in 1:4){
10   plot(0:1,0:1,cex=0,axes=FALSE,main=c("D00","D10","D01","D11")[i])
11   rasterImage(img[[i]],0,0,1,1)
12 }
13 # produce grayscale images
14 get.gray = function(img, w=c(0.22,0.71,0.07)){
15   img[,1]*w[1] + ... # complete this line
16 }
17 img0 = rep(list(NA), 4)
18 for(i in 1:4){
19   img0[[i]] = get.gray(img[[i]])
20 }
21 # plot of grayscale images
22 par(mfrow=c(1,4),mar=c(0,0,1,0))
23 for(i in 1:4){
24   plot(0:1,0:1,cex=0,axes=FALSE,main=c("D00","D10","D01","D11")[i])
25   rasterImage(...) # complete this line
26 }

```

2. (20%) Complete the following R code to produce the empirical CDFs of the grayscale pixel values of the four cases on the same plot.

```

1 col = c("royalblue", "firebrick3", "forestgreen", "goldenrod2")
2 par(mfrow=c(1,1), mar=c(4.5, 5, 3, 2))
3 for(j in 1:4){
4   data = c(img0[[j]])
5   plot(ecdf(...), # complete this line
6        col=col[j], add=j>1, ylab=expression(widehat(F)(x)), xlab=expression(x))
7 }
8 treatment = c("(a) Control", "(b) EB 1089", "(c) IR", "(d) EB 1089 + IR")
9 legend("bottomright", treatment, col=col, lwd=3, bg="white")

```

3. (20%) Consider the cancer treatments (b) EB 1089 and (c) IR. We suspect that the cancer treatments (b) and (c) have different effects on killing the cancer cells. Denote the CDFs of the grayscale pixel values of the treatments (b) and (c) by

$$F_B(\cdot) \quad \text{and} \quad F_C(\cdot).$$

By Sundaram et al. (2003), a lower value of the CDF indicates a lighter image because of the presence of a higher fraction of apoptotic cells. So, we may statistically verify this claim by testing whether the distributions of the grayscale pixel values are the same. Write down the null and alternative hypotheses.

4. (20%) Use your preferred nonparametric method to test the hypotheses stated in part 4.

Hints: See Remark 4.1. Don't read the hints unless you have no ideas and have tried for more than 15 mins.

Exercise 4.2 (★★★ — One-sided distribution test (20%)). Revisit Exercise 4.1. Suppose we suspect the cancer treatment (d) is better than treatment (a). This exercise walks you through the procedures and logic for performing one-sided distribution test. Assume that the CDFs of the grayscale pixel values of the treatments (a) and (d), denoted as $F_A(\cdot)$ and $F_D(\cdot)$, respectively, are continuous and strictly increasing. Let

$$A \sim F_A(\cdot) \quad \text{and} \quad D \sim F_D(\cdot)$$

be independent RVs. Suppose that

$$F_A(t) > F_D(t), \quad t \in \mathbb{R}.$$

1. (⊗) Prove that we can represent A and D as $A = F_A^{-1}(U)$ and $D = F_D^{-1}(V)$, where $U, V \stackrel{\text{iid}}{\sim} \text{Unif}(0, 1)$.
2. (⊗) Prove that $F_A^{-1}(u) < F_D^{-1}(u)$ for all $u \in (0, 1)$.
3. (⊗) Prove that $E(A) < E(D)$.

So, the takeaway of the above properties is that

$$\boxed{F_A(t) > F_D(t), \quad t \in \mathbb{R}} \quad \Rightarrow \quad \boxed{E(A) < E(D)}.$$

Note that the inequality sign is reverse! So, we should test

$$H_0 : F_A(t) = F_D(t) \text{ for all } t \in \mathbb{R} \quad \text{against} \quad H_1 : \sup_{t \in \mathbb{R}} \{F_A(t) - F_D(t)\} > 0,$$

where H_1 essentially means that the image D00.jpg corresponding to treatment (a) is darker than the image D11.jpg corresponding to treatment (d), implying that treatment (d) is better than treatment (a). So,

$$\boxed{\text{Worse treatment}} \Leftrightarrow \boxed{\text{Darker image}} \Leftrightarrow \boxed{\text{Smaller pixel values}} \Leftrightarrow \boxed{\text{Larger CDF values}}.$$

Let $\hat{F}_A(\cdot)$ be the empirical CDF based on the pixel values A_1, \dots, A_m , and $\hat{F}_D(\cdot)$ be the empirical CDF based on the pixel values D_1, \dots, D_n . The one-sided two-sample Kolmogorov–Smirnov (KS) test statistic is

$$K_{m,n}^+ = \sqrt{\frac{mn}{m+n}} \sup_{t \in \mathbb{R}} \{\hat{F}_A(t) - \hat{F}_D(t)\} = \sqrt{\frac{mn}{m+n}} \max_{i=1, \dots, m} \left\{ \frac{i}{m} - \hat{F}_D(A_{(i)}) \right\}, \quad (4.2)$$

where the second equality is true when there is no tied observations. Under H_0 , as $m, n \rightarrow \infty$,

$$K_{m,n}^+ \xrightarrow{d} \sup_{t \in [0,1]} (\mathbb{B}_t - t\mathbb{B}_1) =: K^+,$$

where \mathbb{B}_t is a standard Brownian motion on $[0, 1]$. In R, (4.2) can be computed by the following R-code:

```

1 A = ...      # the pixel value of image corresponding to treatment A
2 D = ...      # the pixel value of image corresponding to treatment D
3 ks.test(A, D, alternative="greater")

```

4. (20%) Find the p -value of the above test. What is your conclusion?

Hints: See Remark 4.2. Don't read the hints unless you have no ideas and have tried for more than 15 mins.

Remark 4.1 (Hints for Exercise 4.1).

3. Check `help(ecdf)` for more information about the function `ecdf`. Example 5.14 is related.
4. The alternative hypothesis is the hypothesis that you suspect true.
5. You may use two-sample two-sided KS test. See Examples 5.10–5.11.

Remark 4.2 (Hints for Exercise 4.2).

1. Try probability inverse transform; see Examples 5.4 and B.25.
2. Apply some inverse function on both side, and then change of variable.
3. If the integrand is smaller, then the integral is ...

It is not hard to write the code for computing the one-sided two-sample KS test statistic. You may try the following comparison:

```

1 > ks.test.statistic.2sided = function(x,y){
2 +   m = length(x)
3 +   n = length(y)
4 +   FY_at_X = ecdf(y)(sort(x))      # ecdf(y) is a function. We can compute, e.g., ecdf(y)(0.4)
5 +   K_plus = sqrt(n/(1+n/m)) * max( (1:m)/m - FY_at_X )
6 +   K_minus = sqrt(n/(1+n/m)) * max( FY_at_X - (0:(m-1))/m )
7 +   K = max(K_plus, K_minus)
8 +   c(two.sided=K, less=K_minus, greater=K_plus)
9 + }
10 >
11 > # simulated dataset
12 > n = 5
13 > m = 8
14 > x = rnorm(n)
15 > y = rnorm(m)
16 >
17 > # self-written version
18 > ks.test.statistic.2sided(x,y)
19 two.sided    less    greater
20 0.5262348 0.3069703 0.5262348
21 >
22 > # built-in version
23 > c(ks.test(x,y, alternative="two.sided")$statistic,
24 +   ks.test(x,y, alternative="less")$statistic,
25 +   ks.test(x,y, alternative="greater")$statistic) * sqrt(m*n/(m+n))
26      D      D^-      D^+
27 0.5262348 0.3069703 0.5262348

```

INSTRUCTIONS: Please follow the instructions stated on the last page of assignment 1.