

Homework on re-sampling methods

P8160 Advanced Statistical Computing

In this homework, we require the use of parallel computing codes for your implementations.

Problem 1: a randomized trial on an eye treatment

An ophthalmologist designed a randomized clinical trial to evaluate a new laser treatment in comparison to the traditional one. The response is visual acuity, measured by the number of letters correctly identified in a standard eye test. 20 patients have both eyes eligible for laser treatment. The ophthalmologist randomized the two laser treatments (new vs traditional) to the two eyes of those patients (i.e. one eye received the new laser treatment and the other receive traditional laser treatment). Another 20 patients had only one suitable eye, so they received one treatment allocated at random. So we have a mixture of paired comparison and two-sample data.

```
> blue <- c(4,69,87,35,39,79,31,79,65,95,68,
            62,70,80,84,79,66,75,59,77,36,86,
            39,85,74,72,69,85,85,72)
> red <-c(62,80,82,83,0,81,28,69,48,90,63,
          77,0,55,83,85,54,72,58,68,88,83,78,
          30,58,45,78,64,87,65)
> acui<-data.frame(str=c(rep(0,20),
                           rep(1,10)),red,blue)
```

Answer the following question:

- (1) The treatment effect of the new laser treatment is defined as

$$E(Y \mid \text{trt} = \text{new}) - E(Y \mid \text{trt} = \text{traditional}).$$

Estimate the treatment effect using the collected data.

- (2) Use bootstrap to construct 95

Problem 2

The Galaxy data consist of the velocities (in km/sec) of 82 galaxies from 6 well-separated conic sections of an unfilled survey of the Corona Borealis region. The structure in the distribution of velocities corresponds to the spatial distribution of galaxies in the far universe. In particular, a multimodal distribution of velocities indicates a strong heterogeneity in the spatial distribution of the galaxies and thus is seen as evidence for the existence of voids and superclusters in the far universe.

Statistically, the question of multimodality can be formulated as a test problem

$$H_0 : n_{\text{mode}} = 1 \quad \text{vs} \quad H_a : n_{\text{mode}} \geq 1$$

where n_{mode} is the number of modes of the density of the velocities.

Considered nonparametric kernel density estimates

$$\hat{f}_{K,h}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

It can be shown that the number of modes in $\hat{f}_{K,h}(x)$ decreases as h increase. Let H_1 be the minimal bandwidth for which $\hat{f}_{K,H_1}(x)$ is unimodal. In the galaxy data, $h_1 = 3.05$

Since multimodal densities need more smoothing to become unimodal, the minimal bandwidth H_1 can be used as a test statistic, and one reject the null hypothesis if

$$\text{Prob}(H_1 > h_1) \leq \alpha$$

To evaluating the distribution of H_1 under the null, one could use the following bootstrap algorithm

1. draw B bootstrap samples if size n from $\hat{f}_{K,h_1}(x)$
2. for each bootstrap, find $h_1^{*(b)}$, the smallest h for which this bootstrap sample has just 1 mode
3. approximate p-value of test is $\frac{\#h_1^{*(b)} > h_1}{B}$

Implement the algorithm above in R, apply it to the galaxy data, and report your findings. You may find the following R codes helpful.

```
library(MASS)
data(galaxies)
plot(density(galaxies/1000, bw=1.5))
plot(density(galaxies/1000, bw=3.5))

#calculate the number of modes in the density
den <- density(galaxies/1000, bw=1.5)
den.s <- smooth.spline(den$x, den$y, all.knots=TRUE, spar=0.8)
s.1 <- predict(den.s, den.s$x, deriv=1)
nmodes <- length(rle(den.sign <- sign(s.1$y))$values)/2
```

Problem 3 (the breast cancer sutdy):

The data *breast-cancer2.csv* have 569 patients. The first column **ID** lables individual breast tissue images; The second column **Diagnosis** identifies if the image is coming from cancer tissue or benign cases (M=malignant, B = benign). There are 357 benign and 212 malignant cases. The other 10 columns correspond to mean of the distributions of the following 10 features computed for the cellnuclei;

- radius (mean of distances from center to points on the perimeter)
- texture (standard deviation of gray-scale values)
- perimeter
- area
- smoothness (local variation in radius lengths)
- compactness ($perimeter^2/area - 1$)
- concavity (severity of concave portions of the contour)
- concave points (number of concave portions of the contour)
- symmetry
- fractal dimension ("coastline approximation" - 1)

Consider a logistic LASSO regression to predict cancer cases based on the image features.

1. Propose and implement a 5-fold cross-validation algorithm to select the turning parameter in the logistic LASSO regression. We call the logistic-LASSO with CV-selected λ as the "optimal" logistic LASSO; The R function for logistic-LASSO
2. Using the selected predictors from the "optimal" logistic LASSO to predict the probability of malignant for each of the images (Note that estimates from logistic-Lasso are biased. You need to re-fit the logistic regression with the selected predictors to estimate the probability.) How well the predictors classify the images?
3. Using the bootstrapping smoothing idea to re-evaluate the probabilities of malignant. How well the new predictors classify the images?
4. Writ a summary of your findings.