```r
 1   # Testing means between groups.
 2
 3   # 1. Comparison between two groups: Two sample t test
 4   # Data are now from two groups, x_1,...,x_n and y_1,...,y_T.
 5   # Assume these two groups follow N(mu1,sigma1^2) and N(mu2,sigma2^2)
 6   # H0: Two groups have the same mean (mu1=mu2)
 7   # t = (x_bar - y_bar)/(SEDM), SEDM = sqrt(sigma1_hat^2/n + sigma2_hat^2/T)
 8
 9   data.set <- read.csv("county_data.csv", stringsAsFactors = FALSE)
10   head(data.set)
11
12   data_1 <- data.set[data.set$State == "California" | data.set$State == "Connecticut",]
13   data_1$State <- factor(data_1$State, levels=c("California", "Connecticut"))
14   boxplot(data_1$Unemployment~data_1$State)
15
16   t.test(data_1$Unemployment[data_1$State=="California"],data_1$Unemployment[data_1$State==
     "Connecticut"], var.equal=FALSE)
17
18   # In textbooks, it is usually assumed that the variances of the two groups are the same.
19
20   t.test(data_1$Unemployment[data_1$State=="California"],data_1$Unemployment[data_1$State==
     "Connecticut"], var.equal=TRUE)
21
22   # 2. Comparison among more than two groups.
23   # Analysis of variance (ANOVA)
24   # Let xgi denote observation no. i in group g.
25   # We can decompose the observations as xgi = x_bar + (xg_bar - x_bar) + (xgi - xg_bar)
26   # xg_bar - x_bar: deviation of group mean from the population mean
27   # xgi - xg_bar: deviation of observation from the group mean
28   # The corresponding model   Xgi = mu + alpha_g +egi,  egi~N(0,sigma^2)
29   # From xgi - x_bar = (xg_bar - x_bar) + (xgi - xg_bar)
30   # Total variation = sum_g sum_i (xgi - x_bar)^2
31   # Within variation = sum_g sum_i (xgi - xg_bar)^2
32   # between variation = sum_g sum_i ng(xg_bar - x_bar)^2
33   # MSw = sum_g sum_i (xgi - xg_bar)^2/(n-G) is an estimate of sigma^2
34   # H0 : all the group means are the same.
35   # If H0 is true, MSb = sum_g sum_i ng(xg_bar - x_bar)^2/(G-1) is also the estimate of
     sigma^2
36
37   # F = MSb / MSw ~F(G-1, n-G)
38
39   data_2 <- data.set[data.set$State %in% c("California", "Connecticut", "Alabama",
     "Ohio"),]
40   data_2$State <- factor(data_2$State, levels=c("Alabama",
     "California","Connecticut","Ohio"))
41   str(data_2)
42
43   boxplot(data_2$Unemployment~data_2$State)
44   anova(lm(data_2$Unemployment~data_2$State))
45
46   # In the outcome: Residual is the within group variation, data_2$State is the between
     group variation.
47
48   # You can also do this test based on the regression coefficients.
49   reg <- lm(data_2$Unemployment~data_2$State)   # Categorial variables (factors) are used
     as dummies.
50   summary(reg)
51
52   # Pairwise comparison. Which pair of states have different means?
53
54   pairwise.t.test(data_2$Unemployment, data_2$State, p.adj="bonferroni")
55   # In multiple testing, use this to be conservative.
56
57   # Exercise
58
59   college <- read.csv("College.csv")
60
61   # College data: Demographic characteristics, tuition, and more for USA colleges.
62   # Private: Public/private indicator
```

```
63    # Apps: Number of applications received
64    # Accept: Number of applicants accepted
65    # Enroll: Number of new students enrolled
66    # Top10perc: New students from top 10 % of high school class
67    # Top25perc: New students from top 25 % of high school class
68    # F.Undergrad: Number of full-time undergraduates
69    # P.Undergrad: Number of part-time undergraduates
70    # Outstate: Out-of-state tuition
71    # Room.Board: Room and board costs
72    # Books: Estimated book costs
73    # Personal: Estimated personal spending
74    # PhD: Percent of faculty with Ph.D.'s
75    # Terminal: Percent of faculty with terminal degree
76    # S.F.Ratio: Student/faculty ratio
77    # perc.alumni: Percent of alumni who donate
78    # Expend: Instructional expenditure per student
79    # Grad.Rate: Graduation rate
80
81    # 1. Compare the distributions of "personal" between private school and public school.
      For this, you can first draw box plots and do 2 sample t test.
82    # 2. Divide the colleges into three groups based on Top10perc. Make the group size to
      be the same with each other.
83    #    Compare the mean of "Grade.Rate" among these four groups. If you conclude there is
      any difference, identify which pair
84    #    of groups have different means.
```