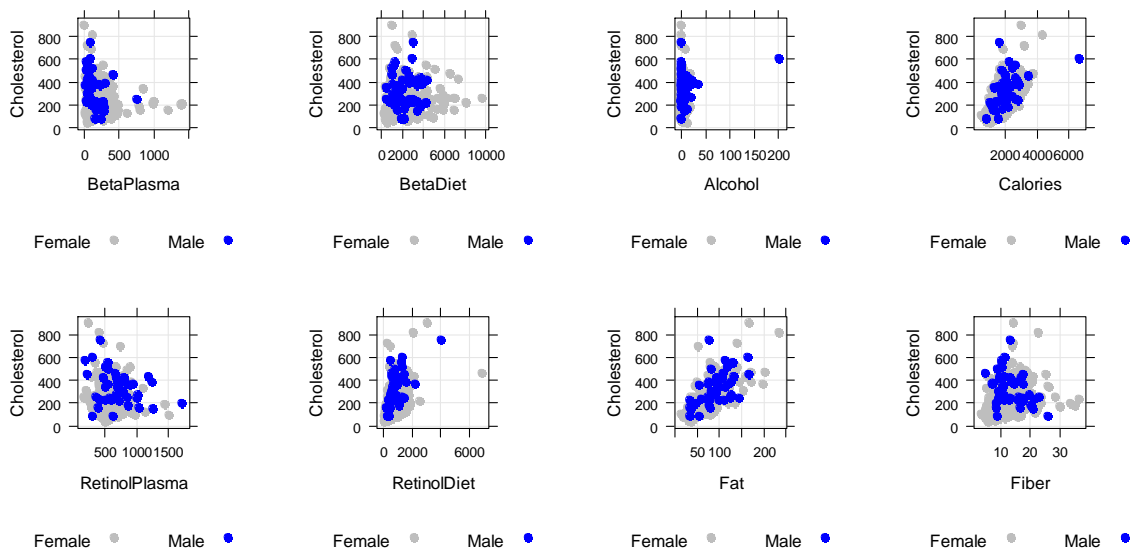# Assignment #6: Fitting and Interpreting ANCOVA Models

**Points: 100 points**

**Data:**  The data for this assignment is the nutrition study data set. This data will be made available by your instructor.
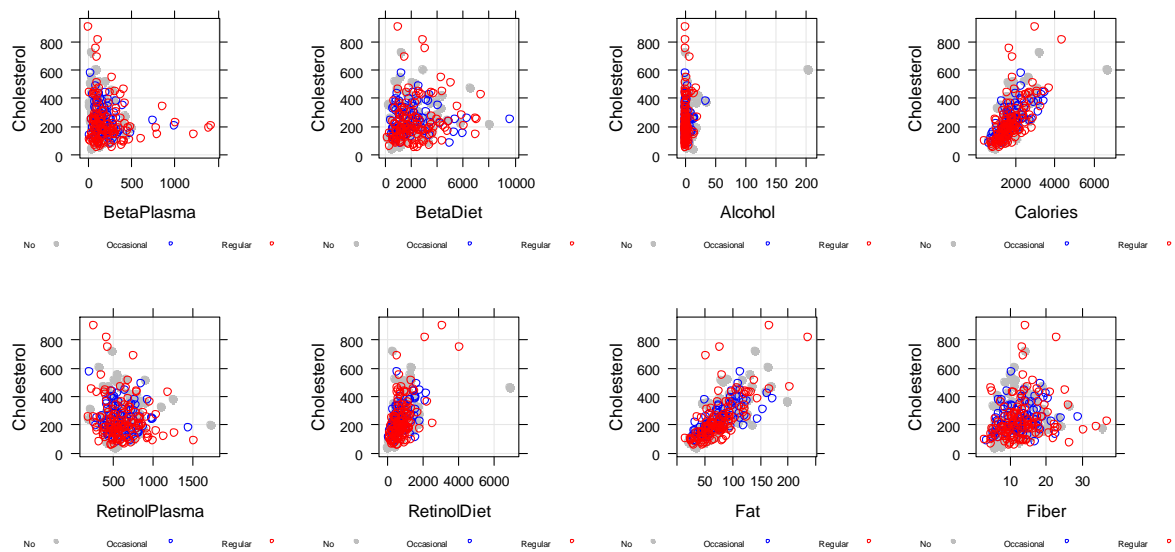
**Assignment Instructions:**

In this assignment we are going to look at a variety of statistical graphics that we may use in the Analysis of Covariance, and then fit some basic ANCOVA models and attempt to understand what they are doing and the output that they produce. We will use a starter script to guide the model fitting so that we all fit the same models correctly, and hence we all get the same output. We will use the output to answer questions in the Assignment #6 template. We will record our answers in green, convert the template to pdf, and submit the pdf as the solution. Hence, this is a question and answer assignment and not a report assignment.
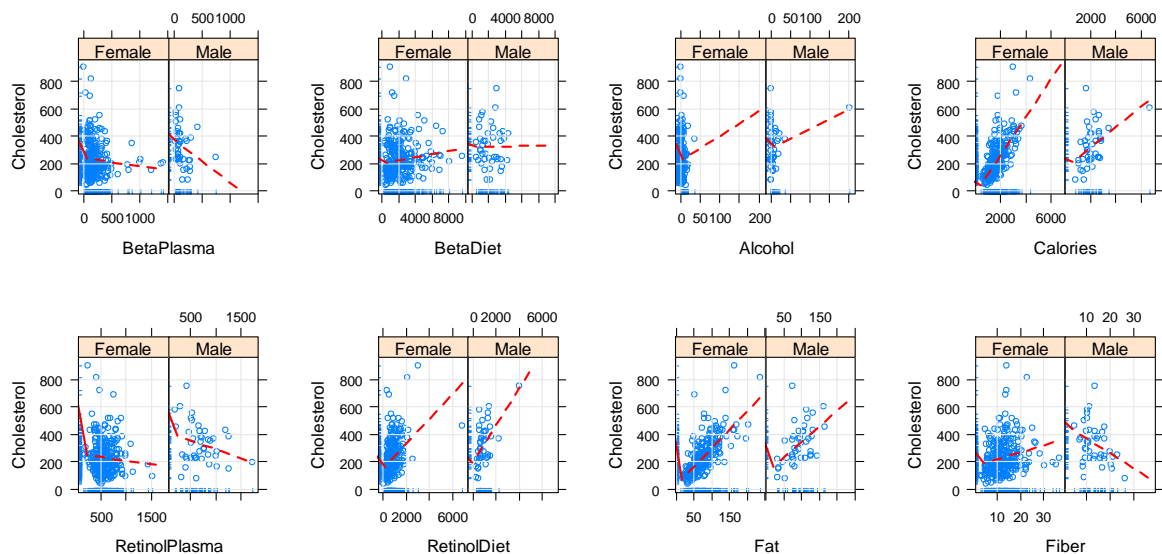


(1) (5 points) If we had to pick one point in all of these plots conditioned on Gender to be an outlier, then which point would we pick? Find the observation in the data set and provide the ID number?

I would pick the Male outlier point on the Alcohol vs Cholesterol chart. This corresponds to ID 62.

(2) (5 points) Do these plots conditioned on vitamin use suggest that vitamin use affects cholesterol? Why or why not?

There doesn't seem to be definite trends that can be identified based on frequency of vitamin use. The three types of points are clustered together and don't indicate a difference based on purely visual analysis.



(3) (10 points) Which two variables appear to have the best regression lines?

(4) (10 points) The variable retinol diet appear to have two regression lines with nice slopes to them. Should we trust those regression lines?  Are there reasons to consider those regression lines to be bad regression lines?

The regression lines for retinol diet seem to follow the basic trend of the variable, but they might not be completely accurate. The female regression line is pulled to a lower slope because of one observation with a high retinol diet value but a mid-cholesterol value. The male regression line is also pulled to a smaller slope because of one of these observations.

(5) (10 points) Are there any plots that do not make any sense?  Hint:  Sometimes we condition on a variable and then the relationship between the predictor variable and the response variable becomes nonsense.  Do we see that behavior in any of the plots?

Some of the variables (BetaDiet, RetinolPlasma) have data points which are spread in a way where a definite pattern is not identifiable. For other plots like Alcohol, the scale makes the relationship between predictor and response very difficult to tell.

We fit ANCOVA models in order to specify regression models with segmented effects on a pooled sample, that is, instead of fitting many separate linear regression models on subsets of the data, we fit a single complex model on the pooled sample.  These specifications allow shared effects and separate effects across the segments.  For the interpretation of ANCOVA models we typically decompose the ANCOVA models into separate models.  Let's fit some ANCOVA models, and see if we can understand what they are.

(6) (10 points) Fit the following ANCOVA model.  Write out the regression equations for Male and Female.

model.1 <- lm(Cholesterol ~ Gender + Fat, data=my.df);

Male: Cholesterol = 76.7355 + 2.6775*Fat

Female: Cholesterol = 29.9715 + 2.6775*Fat

(7) (10 points) Fit the following ANCOVA model.  Write out the regression equations for Male and Female.

model.2 <- lm(Cholesterol ~ Gender:Fat, data=my.df);

Male: Cholesterol = 33.6884 + 3.0383*Fat

Female: Cholesterol = 33.6884 + 2.6465*Fat

(8) (10 points) These two ANCOVA models have different model specifications.  Compute the Mean Absolute Error (MAE) and the Mean Square Error (MSE) for both model.1 and model.2.  Which model should we prefer?

MAE model 1 = 62.11234

MSE model 1 = 8372.765

MAE model 2 = 62.3326

MSE model 2 = 8454.032

Although very close, model 1 has slightly lower error, which would make it the preferred model

(9) (10 points) Fit the following ANCOVA model.  Write out the regression equations for Male and Female.

my.df$Male <- ifelse(my.df$Gender=='Male',1,0);

model.3 <- lm(Cholesterol ~ Male + Fat + Male:Fat, data=my.df);

Male: Cholesterol = 115.9322 + 2.26*Fat
Female: Cholesterol = 25.1472 + 2.7423*Fat

(10) (10 points) Compute the Mean Absolute Error (MAE) and the Mean Square Error (MSE) for model.3. Which of the three models should we prefer?

MAE model 3 = 62.34995
MSE model 3 = 8343.083

This model has about the same MAE as model 1 and 2, but it has a lower MSE than both. This might indicate that model 3 is better.

(11) (10 points) How should we interpret the regression model output for model.3?  Does model.3 suggest that each Gender should have its own intercept and its own slope?

The model is not suggesting each gender gets its own intercept and slope. The default intercept/fat values are for females since we are building our model based on gender. The male estimates are added to the initial coefficients to signify the male component of the regression equation.

Appendix

```
100     col=col.1,
101     pch=pch.1, cex=1, type=c('p','g'),
102     layout=c(1,1), aspect=1.0,
103     key=custom.key
104     )
105
106  plot.8 <- xyplot(Cholesterol ~ Fiber, groups=Gender, data=my.df,
107     col=col.1,
108     pch=pch.1, cex=1, type=c('p','g'),
109     layout=c(1,1), aspect=1.0,
110     key=custom.key
111     )
112
113
114  plot(plot.5, split=c(1,1,2,2))
115  plot(plot.6, split=c(2,1,2,2), newpage=FALSE)
116  plot(plot.7, split=c(1,2,2,2), newpage=FALSE)
117  plot(plot.8, split=c(2,2,2,2), newpage=FALSE)
118
119  my.df$ID[my.df$Alcohol == max(my.df$Alcohol)]
120
121
122 ▾ ############################################################################
123  # Conditional scatter plots - Vitamin Use;
124 ▾ ############################################################################
125
126  pch.1 <- c(19,1,1);
127  col.1 <- c('grey','blue','red');
128
129  custom.key <- list(title='',space='bottom',columns=3,
130     text=list(levels(my.df$VitaminUse)),
131     points=list(pch=pch.1,col=col.1),
132     cex=0.5
133     )
134
135  plot.1 <- xyplot(Cholesterol ~ BetaPlasma, groups=VitaminUse, data=my.df,
136     col=col.1,
137     pch=pch.1, cex=1, type=c('p','g'),
138     layout=c(1,1), aspect=1.0,
139     key=custom.key,
140     )
141
142  plot.2 <- xyplot(Cholesterol ~ BetaDiet, groups=VitaminUse, data=my.df,
143     col=col.1,
144     pch=pch.1, cex=1, type=c('p','g'),
145     layout=c(1,1), aspect=1.0,
146     key=custom.key
147     )
148
149  plot.3 <- xyplot(Cholesterol ~ RetinolPlasma, groups=VitaminUse, data=my.df,
150     col=col.1,
151     pch=pch.1, cex=1, type=c('p','g'),
152     layout=c(1,1), aspect=1.0,
153     key=custom.key
154     )
155
156  plot.4 <- xyplot(Cholesterol ~ RetinolDiet, groups=VitaminUse, data=my.df,
```

```r
280    cex=1, type=c('p','g'),
281    layout=c(2,1), aspect=1.5,
282    panel=my.panel
283    )
284
285  plot.8 <- xyplot(Cholesterol ~ RetinolDiet|Gender, data=my.df,
286    cex=1, type=c('p','g'),
287    layout=c(2,1), aspect=1.5,
288    panel=my.panel
289    )
290
291
292  plot(plot.5, split=c(1,1,2,2))
293  plot(plot.6, split=c(2,1,2,2), newpage=FALSE)
294  plot(plot.7, split=c(1,2,2,2), newpage=FALSE)
295  plot(plot.8, split=c(2,2,2,2), newpage=FALSE)
296
297
298
299
300  ########################################################################
301  # Use ANCOVA to test for a difference in slopes for Gender;
302  ########################################################################
303
304  model.1 <- lm(Cholesterol ~ Gender + Fat, data=my.df);
305  summary(model.1)
306
307
308  model.2 <- lm(Cholesterol ~ Gender:Fat, data=my.df);
309  summary(model.2)
310
311  #MAE and MSE for model 1
312  d = my.df$Cholesterol - predict(model.1)
313  mae1 = mean(abs(d))
314  mse1 = mean((d)^2)
315
316  d2 = my.df$Cholesterol - predict(model.2)
317  mae2 = mean(abs(d2))
318  mse2 = mean((d2)^2)
319
320  my.df$Male <- ifelse(my.df$Gender=='Male',1,0);
321
322  model.3 <- lm(Cholesterol ~ Male + Fat + Male:Fat, data=my.df);
323  summary(model.3)
324
325  d3 = my.df$Cholesterol - predict(model.3)
326  mae3 = mean(abs(d3))
327  mse3 = mean((d3)^2)
328
329
330
331
332
```