## Lecture 5
## Design of Experiments and Confidence Ellipses

In this lecture I'd like to discuss some very general considerations having to do with design of experiments, relating them to the linear regression model and eventually to their implications for making inferences about models with elliptical confidence regions.

In light of these considerations, it is perhaps useful to review some basic facts about confidence regions for parameters in the classical linear regression setting. Suppose that we have a simple linear model with two covariates:

$$y_i = \beta_0 + x_{1i}\beta_1 + x_{2i}\beta_2 + u_i$$

we know that for $u$ spherically normal,

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1})$$

so the variance of any linear contrast $\alpha^\top \hat{\beta}$ is given by evaluating the quadratic form, $\sigma^2 \alpha^\top (X^\top X)^{-1}\alpha$. When $x_1$ and $x_2$ are positively correlated then $\hat{\beta}_1$ and $\hat{\beta}_2$ will be negatively correlated. This implies that we will be able to estimate the sum of the $\beta$'s well, but not their difference.

To illustrate this effect consider the following example from Malinvaud's (1970) classic textbook. We have the following model of French imports:

$$y_t \;=\; \underset{(0.006)}{0.133}x_{1t} \;+\; \underset{(0.110)}{0.550}x_{2t} \;+\; \underset{(0.200)}{2.10}\,x_{4t} \;-\; \underset{(1.27)}{5.92} \tag{1}$$

where $y_t$ is French imports, $x_{1t}$ is gdp, $x_{2t}$ is investment, $x_{3t}$ is consumption, and $x_{4t}$ is dummy variable for EC membership. All variables are in millions of French Francs in 1959 prices. In this model we are able to make reasonably precise estimates of the effect of growth of gdp and investment on imports with 95 percent confidence intervals (respectively)

$$\beta_1 \;\in\; (0.121, 0.145)$$

$$\beta_2 \;\in\; (0.33, 0.77)$$

However, if we introduce the consumption variable $x_{3t}$, we obtain,

$$y_t \;=\; -\underset{(0.051)}{0.021}x_{1t} \;+\; \underset{(0.087)}{0.559}x_{2t} \;+\; \underset{(0.077)}{0.235}x_{3t} \;+\; \underset{(0.16)}{2.10}x_{4t} \;-\; \underset{(1.38)}{9.79} \tag{2}$$

1

But now note that the confidence interval for $\beta_1$ is (-.123,.081). What happened? Roughly speaking, we will see that when, as in this example the independent variables exhibit an approximately linear relationship, here $x_3 \equiv \gamma x_1$, then the "regression" is incapable of precisely estimating the separate effects of the two variables. This is made more explicit if we consider confidence elipses for pairs of coefficients. Without the consumption variable we get a quite precise estimate of the gdp effect, but when we include consumption the situation changes radically – we have a very imprecise estimate of the gdp effect – even the sign of the coefficient is in doubt, and the joint confidence ellipse of the gdp and consumption coefficients is very cigar shaped. Given the orientation of the cigar it is clear that we can quite accurately estimate the effect of circumstances in which gdp and consumption move in the same direction, but we are unable to predict what would happen when they moved in opposite directions. Why?

The intuitive explanation of what went wrong is quite simple. In the more complicated model consumption and gdp are very strongly positively correlated: when gdp (income) goes up there is a natural tendency for personal consumption to also rise. The regression model as specified in (2) would like to distinguish separate effects for these two variables, but the historical experience represented by the data has never seen a period where gdp and consumption deviated substantially from the pattern we have described. When two $x$ variables are very strongly positively correlated like this, it is not surprising that the regression is able to infer very precisely the sum of their two effects, but is unable to precisely infer the difference in their effects. This is basic message of the cigar shape confidence region.

What general conclusions should we draw from this example? First, we can say that it is generally desirable to have explanatory variables ($x$'s) that are uncorrelated with one another. This is already clear from what we have said about making partial residual plots. If $x$'s are uncorrelated then there is no need for the first stage of the PRP procedure; removing the effect of $x_1$ from $x_2$ is superflous if they are already uncorrelated. When experiments are *designed* so that covariates are uncorrelated then we have the luxury of being able to do the least squares analysis one variable at a time. For better or worse, this is rarely an option in economics, and we are usually faced with $x$'s that can be quite strongly correlated. In such cases regression methods do their best to distinguish the separate effects, but often they are not able to do so very precisely. Fortunately, the standard tools for inference, in particular drawing confidence ellipses for pairs of coefficients accurately reflects this imprecision.

The second crucial experimental design consideration is variability of the

**95% Confidence Ellipse**　　　　**95% Confidence Ellipse**



b.investment

0.115　0.125　0.135　0.145

b.gdp
GDP and Investment on Imports

b.consumption

−0.15　−0.05　0.05
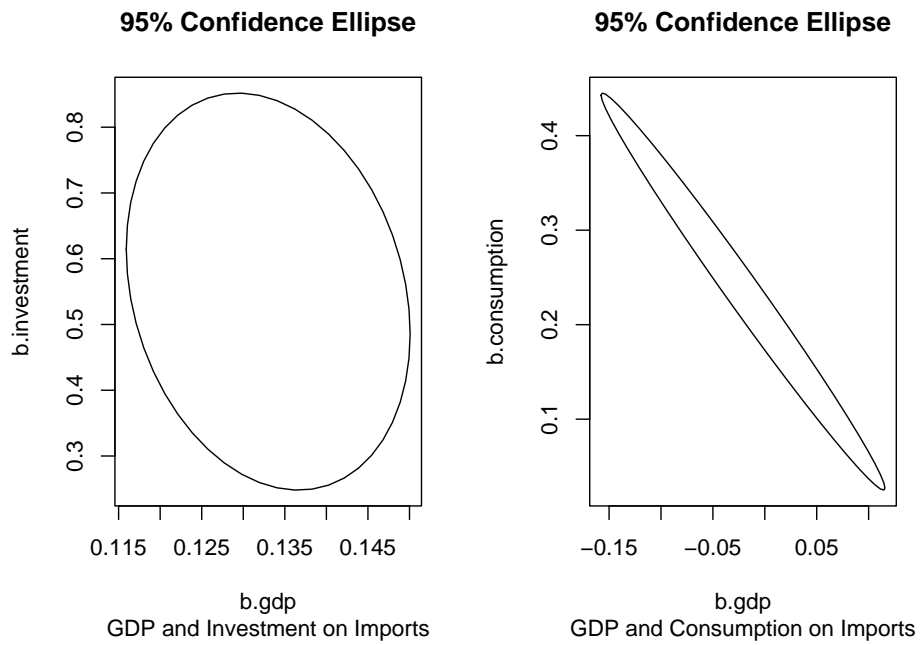
b.gdp
GDP and Consumption on Imports

Figure 1: Confidence ellipses for two pairs of coefficients in the Malinvaud import demand equation. In the left panel the coefficients on gdp and investment are nearly independent, however in the right panel after adding consumption spending, which is quite strongly correlated with gdp, the coefficients of these two variables are very strongly negatively correlated.

covariates. As we saw already in bivariate regression, precision of the slope estimate is inversely proportional to the variance of the observed $x$'s. Thus well designed experiments should try to spread out the $x$ variables as much as possible. Sometimes this is difficult to do because extreme settings of the covariates are difficult or expensive to implement, and in cases where nonlinear effects are suspected there are good reasons to want more uniform spacing, but generally we would like to have more variability in covariates to obtain more precise estimates. Both of these general points are illustrated in the next example.

As a second example consider the problem of jointly estimating confidence intervals for income and price elasticities of gasoline. In Figure 2 we illustrate .90 and .99 confidence ellipses for two estimated gasoline models. One is based on data from 1947-72 prior to the first oil shock, and the other is based on the entire period 1947-88. Several things are evident from the figure. First, the full data set yeilds much more precise estimates (smaller confidence regions). This is to be expected when there is more data, and more especially when there is more variability in the $x$ variables, as was kindly provided by OPEC. Second, the orientation of the ellipses for the full sample is somewhat more aligned with the coordinate axes indicating that there is less correlation between the two elasticities than in the earlier period. This reflects more independent movement of prices in the OPEC period, whereas price and income were more strongly positively correlated in the earlier pre-OPEC period. Finally, and most disturbingly note that the evidence provided by the earlier period is wildly overconfident about precision of the elasticity estimates. While admitting that the price elasticity might be negative, it rules out very strongly the possibility that it could be as small as -0.50, the value obtained using the full data set. Similarly, the confidence in the lower estimate of the income elasticity is also misplaced.

Finally, to conclude this digression, let's consider the relationship between the confidence ellipses that we have seen and the conventional one dimensional confidence intervals. To fix ideas let's consider the simplest possible case: a situation in which we have a two dimensional parameter $\beta$ that happens to be standard normal, i.e. $\beta \sim \mathcal{N}(0, I_2)$. This is a totally artificial situation in which we imagine that $\hat{\beta}$ happens to take the value $(0,0)^\top$ and have covariance matrix, $I_2$. Then we have that

$$P(\beta_1^2 + \beta_2^2 < 5.99) = .95$$

since the sum of squares of the $\beta$'s is $\chi_2^2$. Thus, we get circular confidence regions and the radius of the .95 region is 2.45. Compare the aria of this
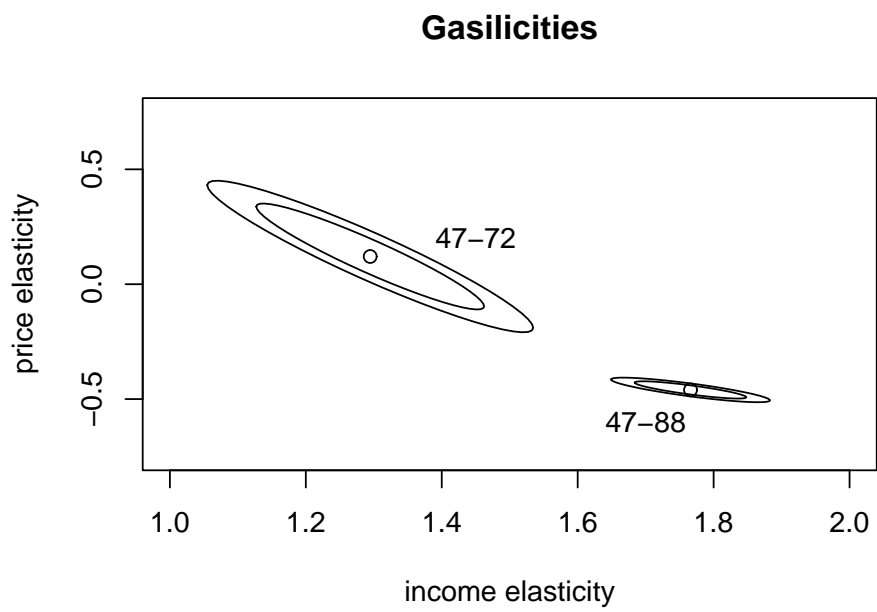
**Gasilicities**



Figure 2: Confidence ellipses for income and price elasticities of gasoline in the U.S.

circle: $\pi r^2 = 18.81$ with the area of the square formed by two .95 confidence intervals for the separate parameters: which has area $(2 \cdot 1.96)^2 = 15.36$. Why is this square smaller than the circle? Hint: Find the the square that contains probability .95 and compare its area with those already computed.