

Chapter 6: Statistical Inference for Regression

6.1 Simple Regression Model

Population regression model:

$$Y_i = \alpha + \beta X_i + \varepsilon_i, \quad i = 1, \dots, n$$

$$y = \alpha + \beta x + \varepsilon$$

$$y = \mathbf{X}\boldsymbol{\beta} + \varepsilon$$

The errors ε_i represent

- omitted explanatory variables
- measurement error
- “inherently random component of Y”

Here we address key assumptions in order to do inference (confidence intervals and tests). But first more on ...

Notation: From now on the text notation is mostly consistent in using capital letters for random variables and lowercase letters to denote observations, as in the first line above (I wrote it with a small “y” before). However, it is still somewhat variable as the regression model is variously written with capital and lowercase x on the right-hand side. A footnote on p. 100 explains that Fox uses lowercase x to stress that the value of x_i is fixed, either literally in an experimental design, or by conditioning on the observed value x_i of X_i . So you see that the explanatory variable could be considered as a random variable or fixed numbers. In addition, chap 9 chooses to denote a vector of random variables by a lower case bold \mathbf{y} as shown above, rather than a bold capital letter.

I tend to use the “hat” notation for estimates where as Fox uses latin letters for estimates of (greek letter) parameters.

Key assumptions in order to do inference (confidence intervals and tests):

- Linearity: $E(y_i) = \alpha + \beta x_i$,
or $E(\varepsilon_i) = E(\varepsilon_i | x_i) = 0$
- Constant variance: $\text{Var}(\varepsilon_i | x_i) = E(Y_i - \alpha - \beta x_i)^2 = \sigma_\varepsilon^2$
- Normality: $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$
 - Question: Do we assume that the Y_i are normally distributed?
- Independence: random errors ε_i and ε_j are statistically (probabilistically) independent. This assumption is generally justified by knowledge of how the data were sampled/collected. For example, observations which represent
 - Random sample of subjects from large population, vs.
 - Time series, vs
 - Sample of children from a school with multiple classrooms. But here, the assumption will depend on what you are measuring. An analysis of `height ~ age`, independence between children may be a reasonable assumption, for for an analysis of `spelling ~ age`, probably not.
- X's are either
 - Fixed by design, or
 - Sampled/measured with error independent of ε_i

Examples for discussion:

- Average income -> Prestige for Canadian occupations
- Education -> Hourly wage rate for Canadian employees
- GDP per capita -> Infant mortality in the UN data base

6.1.2 Properties of least squares estimates:

$$\hat{\alpha} = \bar{Y} - B\bar{X}$$

$$\hat{\beta} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

or, in matrix notation

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

(1) L.S. estimates are linear in Y

$$\hat{\beta} = \sum m_i Y_i, \quad m_i = \frac{x_i - \bar{x}}{\sum (x_j - \bar{x})^2}$$
$$\hat{\alpha} = ? \text{ (exercise)}$$

(2) Sampling variance is

$$\begin{aligned}
\text{Var}(\hat{\beta}) &= \text{Var}\left(\sum m_i Y_i\right) = \sum m_i^2 \text{Var}(Y_i) + 2 \sum_{i < j} m_i m_j \text{Cov}(Y_i, Y_j) \\
&= \frac{1}{\left(\sum (x_j - \bar{x})^2\right)^2} \sum (x_i - \bar{x})^2 \sigma_\varepsilon^2 \\
&= \frac{\sigma_\varepsilon^2}{\sum (x_j - \bar{x})^2} = \frac{\sigma_\varepsilon^2}{(n-1)s_x^2}
\end{aligned}$$

Question: If you are designing a study, how can you make the uncertainty (standard error) of $\hat{\beta}$ small?

Similarly, you can show

$$\begin{aligned}
\text{Var}(\hat{\alpha}) &= \text{Var}\left(\sum q_i Y_i\right) \\
&= \frac{\sigma_\varepsilon^2 \sum x_i^2}{n \sum (x_j - \bar{x})^2}
\end{aligned}$$

which is small(est) when the x's are centered at zero.

(3) If $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$, then

$$\begin{aligned}
\hat{\alpha} &\sim N\left(\alpha, \frac{\sigma_\varepsilon^2 \sum x_i^2}{n \sum (x_j - \bar{x})^2}\right) \\
\hat{\beta} &\sim N\left(\beta, \frac{\sigma_\varepsilon^2}{\sum (x_j - \bar{x})^2}\right)
\end{aligned}$$

But what if the errors cannot be assumed normally distributed?

Recall that $\hat{\beta} = \sum m_i Y_i$, so that by the Central Limit Theorem $\hat{\beta}$ should be approximately large

- (4) By the Gauss-Markov Theorem, of all linear estimators $\beta^* = \sum c_i Y_i$ that are unbiased, $E\beta^* = \beta$, $\hat{\beta}_{LS}$ has the minimum variance --- it is the “most efficient” estimator.
- (5) Under all the assumptions stated above, which can be combined into one statement,

$$\varepsilon_i \underset{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

$(\hat{\alpha}, \hat{\beta})$ are maximum likelihood estimates. (Consider reading Appendix D (online) to the text by Fox, although this section goes into more detail than you need)

6.1.3 Confidence intervals and hypothesis tests:

$\text{Var}(\hat{\beta}) = \frac{\sigma_\varepsilon^2}{(n-1)s_x^2}$ but σ_ε^2 is unknown, so we substitute the estimate s_E^2 and write

$$\hat{\text{Var}}(\hat{\beta}) = \frac{s_E^2}{(n-1)s_x^2} \quad \text{and} \quad \text{s.e.}(\hat{\beta}) = \sqrt{\frac{s_E^2}{(n-1)s_x^2}}$$

Defn: The standard error of a statistics is an *estimate* of the *standard deviation of its sampling distribution*.

Question: What is the definition of a sampling distribution?

*** R simulation of sampling distributions ***

Confidence intervals: (review?)

$$\frac{\hat{\beta} - \beta}{\sqrt{\sigma_\varepsilon^2 / \sum (x_i - \bar{x})^2}} \sim ?$$

$$\frac{\hat{\beta} - \beta}{\sqrt{\text{s.e.}(\hat{\beta})}} \sim t_{(n-2)}$$

$$\Leftrightarrow P\left(-t_{n-2, \alpha/2} \leq \frac{\hat{\beta} - \beta}{\sqrt{\text{se}(\hat{\beta})}} \leq t_{n-2, \alpha/2}\right) = 1 - \alpha \quad (\text{where } t_{n-2, \alpha/2} \text{ is the } (1 - \alpha/2) \text{ \%ile})$$

$$\Leftrightarrow P\left(\hat{\beta} - t_{n-2, \alpha/2} \text{se}(\hat{\beta}) \leq \beta \leq \hat{\beta} + t_{n-2, \alpha/2} \text{se}(\hat{\beta})\right)$$

Same type of calculation applies, of course, for $\hat{\alpha}$.

Hypothesis test: (review?)

$H_0 : \beta = \beta_0$ (most often considering $\beta_0 = 0$)

$$t_0 = \frac{\hat{\beta} - \beta_0}{\text{se}(\hat{\beta})} \sim t_{(n-2)} \text{ under } H_0$$

$$P_{H_0} \left(\left| \frac{\hat{\beta} - \beta_0}{\text{se}(\hat{\beta})} \right| > t_{(n-2), \alpha/2} \right) = \alpha$$

\Rightarrow reject if $t_0 > t_{(n-2), \alpha/2}$

Example: Consider the Davis weight-reported weight data. Test $H_0 : \beta = 1$

Question: How would you test the null hypothesis $H_0 : \beta = 1$ and $\alpha = 0$?

(there are two computational approaches that provide the same answer)

6.2 Multiple Regression

The mathematics becomes a little more interesting for multiple regression problems. (Sects 6.2, 9.3).

Let's first go back and revisit the least squares solution of chapter 5 in a little more detail.

$$\begin{aligned}RSS &= \|\boldsymbol{\varepsilon}\|^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\&= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}\end{aligned}$$

Results for differentiation of a matrix expression:

$$\begin{aligned}\frac{\partial(\mathbf{c}^T \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= \begin{pmatrix} \partial(\mathbf{c}^T \boldsymbol{\beta}) / \partial \beta_1 \\ \vdots \\ \partial(\mathbf{c}^T \boldsymbol{\beta}) / \partial \beta_{k+1} \end{pmatrix} = \begin{pmatrix} c_1 \\ \vdots \\ c_{k+1} \end{pmatrix} = \mathbf{c} \\ \frac{\partial(\boldsymbol{\beta}^T \mathbf{C} \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= 2\mathbf{C}\boldsymbol{\beta}\end{aligned}$$

which leads to

$$\frac{\partial RSS}{\partial \boldsymbol{\beta}} = 2\mathbf{X}^T \mathbf{y} - 2\mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = 0$$

or

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

This works if $(\mathbf{X}^T \mathbf{X})$ has an inverse, which means that \mathbf{X} must have “full column rank”. That is, the columns of \mathbf{X} must be *linearly independent*: no column can be an exact linear combination of the other columns. We’ll see shortly that there are consequences for columns which are nearly linearly dependent.

Next consider the basic properties in matrix notation:

$$\begin{aligned} E(\hat{\beta}) &= ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\mathbf{y}) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta = \beta \end{aligned}$$

$$\text{Var}(\hat{\beta}) = \text{Var}(\mathbf{M}\mathbf{y}) = \mathbf{M} \text{Var}(\mathbf{y}) \mathbf{M}^T$$

This is a matrix of variances and covariances. We are simplifying notation setting $\mathbf{M} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}$.

Because $\text{Var}(\mathbf{y}) = \sigma_{\varepsilon}^2 \mathbf{I}_n$, where \mathbf{I}_n is the $n \times n$ identity matrix, the expression above simplifies to

$$\text{Var}(\hat{\beta}) = \sigma_{\varepsilon}^2 \mathbf{M} \mathbf{M}^T = \sigma_{\varepsilon}^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} = \sigma_{\varepsilon}^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

You should know and be comfortable with the derivation of this fundamental expression.

Under normality, $\hat{\beta} \sim N_{k+1}(\beta, \sigma_{\varepsilon}^2 (\mathbf{X}^T \mathbf{X})^{-1})$.

The standard errors of the estimated regression coefficients are the square roots of the diagonal elements of $\sigma_{\varepsilon}^2 (\mathbf{X}^T \mathbf{X})^{-1}$.

There is a very useful expression for the $\text{Var}(\hat{\beta}_j)$, the j th diagonal element of $\sigma_\varepsilon^2(\mathbf{X}^T\mathbf{X})^{-1}$ that is reported as eqn (6.2). It can be shown (but you will not be required to show) that

$$\begin{aligned}\text{Var}(\hat{\beta}_j) &= \frac{1}{\sqrt{1-R_j^2}} \frac{\sigma_\varepsilon^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_i)^2} \\ &= \frac{\sigma_\varepsilon^2}{\sum_{i=1}^n (x_{ij} - \hat{x}_{ij})^2}\end{aligned}$$

where R_j^2 is the *squared multiple correlation coefficient* for regression of \mathbf{x}_j , the j th column of \mathbf{X} , on all the other columns of \mathbf{X} , and \hat{x}_{ij} represents the fitted values from this auxiliary regression.

The term $\frac{1}{\sqrt{1-R_j^2}}$ is called the *variance-inflation factor*. If this factor is large, meaning that the j th column of \mathbf{X} is well-predicted by (highly correlated with) the other columns of \mathbf{X} , the uncertainty in $\hat{\beta}_j$ will be large. The second line of the expression for $\text{Var}(\hat{\beta}_j)$ tells us the same thing: if the j th column of \mathbf{X} is well-predicted by the other columns of \mathbf{X} , the *conditional* variation in \mathbf{x}_j given its prediction $\hat{\mathbf{x}}_j$ by the other columns, $\sum_{i=1}^n (x_{ij} - \hat{x}_{ij})^2$, will be small and, hence, $\text{Var}(\hat{\beta}_j)$ will be large. *That is*, we will have a problem if the columns of \mathbf{X} are nearly linearly dependent.

Tests of $H_0 : \beta_j = 0$ and confidence intervals for β_j are computed exactly as in simple linear regression, but using the standard errors of $\hat{\beta}_j$ presented above.

Table 9.1 Comparison between simple regression using scalars and multiple regression using matrices

	Simple Linear Regression	Multiple Linear Regression
Model	$Y = \alpha + \beta x + \varepsilon$	$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \varepsilon$
Least-squares estimator	$\hat{\beta} = \frac{\sum x^* Y^*}{\sum (x^*)^2}$	$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
Sampling variance	$V(\hat{\beta}) = \frac{\sigma_\varepsilon^2}{\sum (x^*)^2} = \sigma_\varepsilon^2 (\sum (x^*)^2)^{-1}$	$V(\hat{\boldsymbol{\beta}}) = \sigma_\varepsilon^2 (\mathbf{X}^T \mathbf{X})^{-1}$
Distribution	$\hat{\beta} \sim N\left(\beta, \sigma_\varepsilon^2 (\sum (x^*)^2)^{-1}\right)$	$\hat{\boldsymbol{\beta}} \sim N_{k+1}\left(\boldsymbol{\beta}, \sigma_\varepsilon^2 (\mathbf{X}^T \mathbf{X})^{-1}\right)$

Interpretation of the multiple regression model $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$

- In a *simple linear regression* (one explanatory variable), the coefficient β is called a marginal regression coefficient
- In a *multiple linear regression*, the coefficient β_j is called the partial regression coefficient—the “effect” on the response variable of a one unit change in the corresponding explanatory variable holding constant the values of all other explanatory variables. *[Does this always make sense?]*
- The least squares coefficient represents the average change in Y associated with a one unit change in X_j when all the other X 's are held constant.

R example: Prestige ~ Income + Education + . . .

The ANOVA for multiple regression. As shown in chapter 5

$$\text{TSS} = \text{RegSS} + \text{RSS}$$

where

$$\begin{aligned}\text{TSS} &= \sum (Y_i - \bar{Y})^2 \\ \text{RegSS} &= \sum (\hat{Y}_i - \bar{Y})^2 \\ \text{RSS} &= \sum (Y_i - \hat{Y}_i)^2\end{aligned}$$

The conventional table is:

Source	SS	Df	MS	F
Regression	RegSS	k	RegMS=RegSS/k	RegMS/RMS
Residual	RSS	n-k-1	RMS=RSS/(n-k-1)	
Total	TSS	n-1		

Under $H_0 : \beta_1 = \beta_2 = \dots = \beta_k$, $F = \frac{\text{RegMS}}{\text{RSS}} = \frac{\text{"MSReg"}}{\text{"MSE"}} \sim F_{k, n-k-1}$

Why? The F statistic is the ratio of two independent mean squares (they are independent because they derive from squared lengths of orthogonal vectors in a geometric representation) divided by their degrees of freedom:

$$F = \frac{\chi_k^2 / k}{\chi_{n-k-1}^2 / (n-k-1)}$$

In general, $E(\text{RegMS}) = \sigma_\varepsilon^2 + f(\beta's)$, and under $H_0 : \beta_1 = \beta_2 = \dots = \beta_k$, $E(\text{RegMS}) = \sigma_\varepsilon^2$

We always have $E(\text{RMS}) = \sigma_\varepsilon^2$, so under H_0 we have the ratio of two independent chi-square random variables, both estimating the same thing, and this is the definition of a *central* F-distribution.

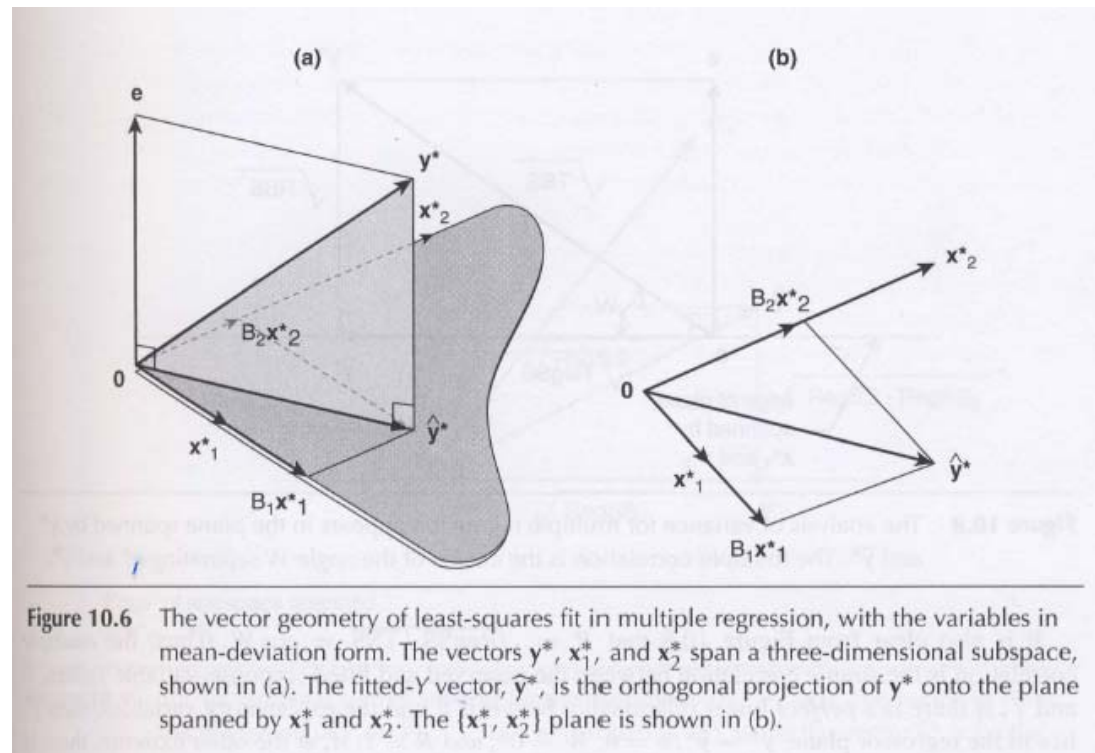
A little algebra shows that we can also express the F-statistic as

$$F = \frac{n-k-1}{k} \times \frac{R^2}{1-R^2}$$

It really helps to see this in matrix notation: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$

In Chap 11 you'll see the notation $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ and the residuals are $(\mathbf{y} - \hat{\mathbf{y}}) = (\mathbf{I}_n - \mathbf{H})\mathbf{y}$, where $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is commonly referred to as the “hat matrix”. It is a projection matrix: $\mathbf{H}\mathbf{H} = \mathbf{H}$, and similarly $(\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H}) = (\mathbf{I} - \mathbf{H})$. And, $\mathbf{H}(\mathbf{I} - \mathbf{H}) = \mathbf{0}$, so these two projection matrices are orthogonal.

This is illustrated in Fig 10.6 for a regression with $k=2$ explanatory variables and the problem represented in mean-centered form. $\hat{\mathbf{y}}^* = \mathbf{H}\mathbf{y}^*$, where \mathbf{H} is the projection matrix defined from the design matrix $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^* & \mathbf{x}_2^* \end{bmatrix}$.



(Dropping the notation for the centered variables)

$$\begin{aligned}
 \mathbf{y}^T \mathbf{y} &= (\hat{\mathbf{y}} + (\mathbf{y} - \hat{\mathbf{y}}))^T (\hat{\mathbf{y}} + (\mathbf{y} - \hat{\mathbf{y}})) \\
 &= (\mathbf{H}\mathbf{y} + (\mathbf{I} - \mathbf{H})\mathbf{y})^T (\mathbf{H}\mathbf{y} + (\mathbf{I} - \mathbf{H})\mathbf{y}) \\
 &= \mathbf{y}^T \mathbf{H}\mathbf{y} + \mathbf{y}^T (\mathbf{I} - \mathbf{H})\mathbf{y} = \text{RegSS} + \text{RSS}
 \end{aligned}$$

Note that the Regression SS can also be written

$$\mathbf{y}^T \mathbf{H}\mathbf{y} = \mathbf{y}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y}$$

The ANOVA table can therefore be written

Source	SS	Df	MS
Regression	$\mathbf{y}^T \mathbf{H}\mathbf{y} = \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y}$	k+1	RegMS=RegSS/(k+1)
Residual	$\mathbf{y}^T (\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{y}^T \mathbf{y} - \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y}$	n-k	RMS=RSS/(n-k)
Total (“uncorrected”)	$\mathbf{y}^T \mathbf{y}$	n-1	

If we write everything in mean-centered form, the same expressions hold (with *’s on all the vectors and matrices to denote mean-centering) with k and (n-k-1) df. The “Total” sum of squares is then “corrected for the mean”,

$$\mathbf{y}^{*T} \mathbf{y}^* = \sum (y_i - \bar{y})^2.$$

The *omnibus* F-test derived from this ANOVA is a test of the significance of the multiple regression, rejected if there is evidence that at least one of the β 's is significantly different from zero.

A modest generalization of this provides a basis for testing a subset of the slopes. In the general framework we define

RSS_1 and $RegSS_1$ the regression sum of squares for the full model

RSS_0 and $RegSS_0$ the residual and regression sums of squares for the reduced model

and

$$F = \frac{(RegSS_1 - RegSS_0) / q}{RSS_1 / (n - k - 1)} = \frac{(RSS_0 - RSS_1) / q}{RSS_1 / (n - k - 1)}$$

In matrix notation,

$$\mathbf{y} = \alpha + \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

where $\boldsymbol{\beta}_1 = (\beta_1, \dots, \beta_q)^T$, $\boldsymbol{\beta}_2 = (\beta_{q+1}, \dots, \beta_k)^T$ and $\mathbf{X} = [\mathbf{X}_1 \mathbf{X}_2]$ with \mathbf{X}_1 and \mathbf{X}_2 having q and $(k - q)$ columns, respectively.

Let

$RegSS_0 = SS(\hat{\boldsymbol{\beta}}_2) = \hat{\boldsymbol{\beta}}_2^T \mathbf{X}_2^T \mathbf{y}$, the regression sum of squares for the *reduced model*.

$RegSS_1 = SS(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\beta}}_2) = \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y}$, the regression sum of squares for the *full model*.

Under $H_0 : \beta_1 = 0$, $F = \frac{(SS(\hat{\beta}_1, \hat{\beta}_2) - SS(\hat{\beta}_2)) / q}{RSS / (n - k - 1)} \sim F_{q, (n-k-1)}$

***If \mathbf{X}_1 and \mathbf{X}_2 are orthogonal, then**

$$SS(\hat{\beta}_1, \hat{\beta}_2) = \hat{\beta}_1^T \mathbf{X}_1^T \mathbf{y} + \hat{\beta}_2^T \mathbf{X}_2^T \mathbf{y}$$

Source	SS	Df	MS
$\mathbf{X}_2 : SS(\hat{\beta}_2)$	$\hat{\beta}_2^T \mathbf{X}_2^T \mathbf{y}$	$k - q$	RegMS ₂
$\mathbf{X}_1 \mathbf{X}_2 : SS(\hat{\beta}_1 \hat{\beta}_2)$	$\hat{\beta}^T \mathbf{X}^T \mathbf{y} - \hat{\beta}_2^T \mathbf{X}_2^T \mathbf{y}$	q	RegMS _{1 2}
Residual	$\mathbf{y}^T (\mathbf{I} - \mathbf{H}) \mathbf{y} = \mathbf{y}^T \mathbf{y} - \hat{\beta}^T \mathbf{X}^T \mathbf{y}$	n-k	RMS=RSS/(n-k)
Total	$\mathbf{y}^{*T} \mathbf{y}^*$	n-1	

* A special case for testing $q = 1$ coefficient, $H_0 : \beta_1 = 0$. F-test or t-test?

$$F_{1, (n-k-1)} = t_{(n-k-1)}^2. \quad \text{Because } F = \frac{n-k-1}{k} \times \frac{R^2}{1-R^2}, \quad t = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}}.$$

Return to R code and anova() for multiple regression.

6.3 Empirical vs Structural Relations

Empirical: Descriptive relationship among variables

Structural: Descriptions from which we intend to infer *causation*, a model of how response scores are actually determined.

$$(a) \quad Y = \alpha' + \beta_1' X_1 + \varepsilon'$$

$$(b) \quad Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

If the relationship between Y and X_1 is well described by the simple linear relationship (a), the fact that $\beta_1' \neq \beta_1$, so that $\hat{\beta}_1'$ may be considered a biased estimate of β_1 , is not necessarily an issue. However, if we really believed that the data were generated by model (b), this bias would be important.

If (b) is true, but we fit (a), then $\varepsilon' = (\beta_2 X_2 + \varepsilon)$, so that X_1 and ε' are correlated if X_1 and X_2 are correlated.

Correlation between the error term and the explanatory variable is a problem leading to bias of $\hat{\beta}_1'$ with respect to the value of β_1 in (b).

On pp. 111-112 Fox derives an expression for β_1 in model (b)

Take expected values of (b), giving

- $\mu_y = \alpha + \beta_1 \mu_1 + \beta_2 \mu_2 + 0$

and subtract this from (b), giving

- $(Y - \mu_y) = \beta_1 (X_1 - \mu_1) + \beta_2 (X_2 - \mu_2) + \varepsilon$

Multiply by $(X_1 - \mu_1)$, giving

- $(X_1 - \mu_1)(Y - \mu_y) = \beta_1 (X_1 - \mu_1)^2 + \beta_2 (X_1 - \mu_1)(X_2 - \mu_2) + (X_1 - \mu_1)\varepsilon$

Take expected value of both sides of this equation:

- $\sigma_{1Y} = \beta_1 \sigma_1^2 + \beta_2 \sigma_{12}$

Thus, solving for β_1 ,

- $\beta_1 = \frac{\sigma_{1Y}}{\sigma_1^2} - \beta_2 \frac{\sigma_{12}}{\sigma_1^2}$

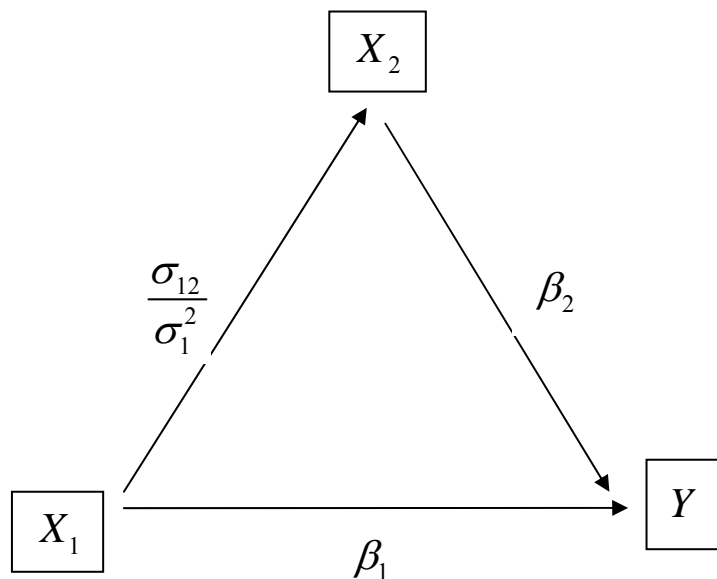
The least squares estimate from model (a) is $\hat{\beta}'_1 = \frac{s_{1Y}}{s_1^2}$, so that $\hat{\beta}'_1$ is an estimate of $\frac{\sigma_{1Y}}{\sigma_1^2} = \beta'_1$, meaning it is biased as

an estimator of the coefficient for model (b), $\beta_1 = \frac{\sigma_{1Y}}{\sigma_1^2} - \beta_2 \frac{\sigma_{12}}{\sigma_1^2}$. i.e., $E(\hat{\beta}'_1) = \beta_1 + \beta_2 \frac{\sigma_{12}}{\sigma_1^2}$.

We have bias if

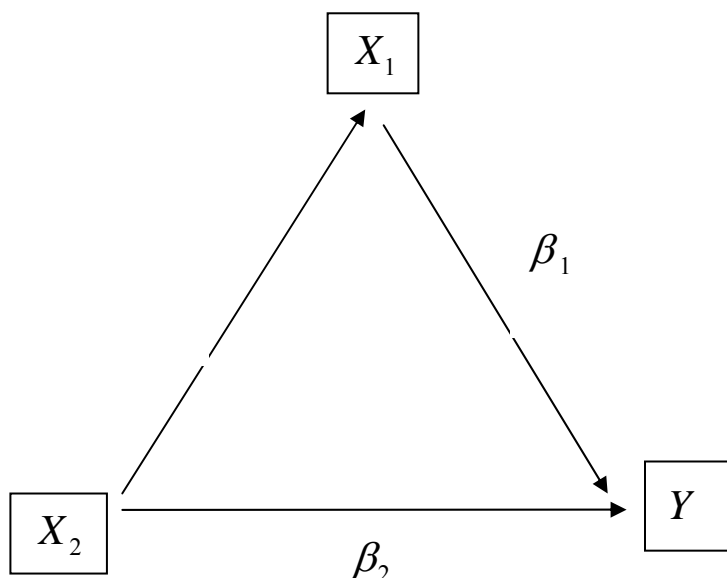
1. X_2 is “relevant” (i.e., $\beta_2 \neq 0$), and
2. X_1 and X_2 are correlated.

One final subtlety of interpretation: the “bias” $\beta_2 \frac{\sigma_{12}}{\sigma_1^2}$ depends on the nature of the causal relationship between X_1 and X_2 . In Fig 6.2(a) $\beta_2 \frac{\sigma_{12}}{\sigma_1^2}$ is seen to be the indirect effect of X_1 on Y through X_2 .



$\frac{\sigma_{12}}{\sigma_1^2}$ is the population slope for the regression of X_2 on X_1 . X_2 is an intervening variable on a causal path through X_2 .

By contrast, in Fig 6.2(b), X_2 is a common prior cause of both X_1 and Y .



In this case the bias is considered a spurious (non-causal) component of the association between X_1 and X_2 .

And in this case it is critical to include the variable X_2 in the analysis— “to control for X_2 ” in examining the association between X_1 and X_2 .

X_2 is sometimes described as a lurking variable. These are what make causal inference from observational studies so difficult.

Omission of X_2 is not necessarily (always) an issue in the case of Fig 6.2(a), depending on whether it is important to understand the direct and indirect effects of X_1 .

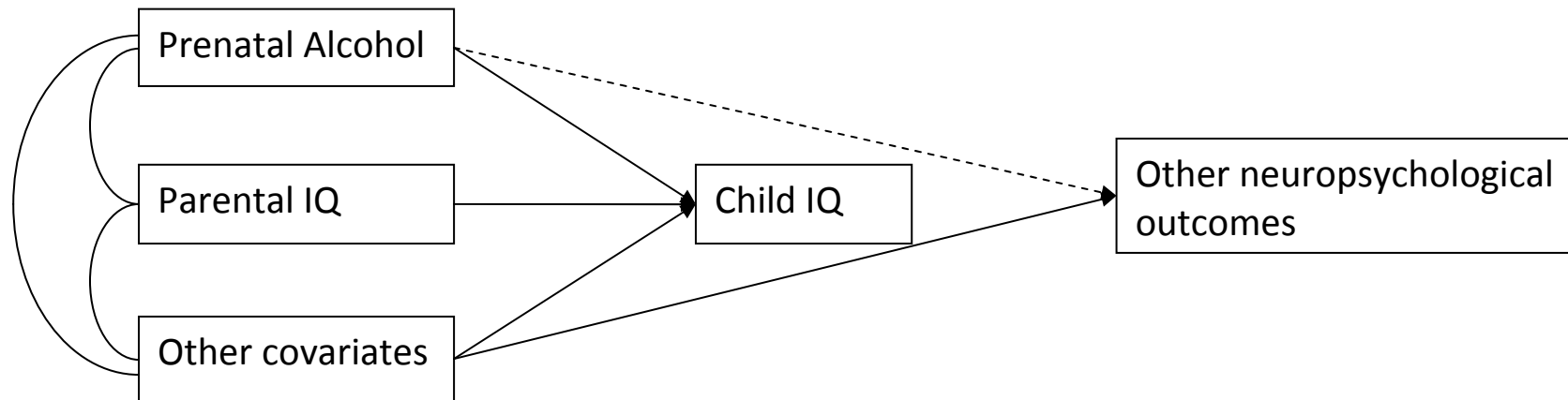
So, you have to think hard about your problem if the scientific question(s) target the effect of a particular explanatory variable “ X_1 ”.

Examples:

- X_1 is a measure of Chinese investment in the U.S. Y is a measure of U.S. stock prices. “ X_2 ”?
- X_1 is an indicator distinguishing two ethnic groups. Y is a measure of public health outcome. “ X_2 ”?

See Fox Exercise 6.8 in addition to Exercise 6.9 assigned for Homework #4.

Another example from my applied work: The effects of Prenatal Alcohol exposure on offspring IQ.



We can't say that Parental IQ causes Maternal drinking (prenatal alcohol exposure), but it is correlated with prenatal alcohol (in a manner that depends on how you measure prenatal alcohol), so it contributes to the observed association of Prenatal Alcohol with Child IQ.

⇒ you must have Parental IQ in the model or else you will have a biased estimate of the effect of Prenatal Alcohol on Child IQ.

What do you do if you haven't measured Parental IQ?

6.4 Measurement Error in Explanatory Variables

The basic multiple regression models we work with assume that the explanatory variables are known without error. Sometimes this is a reasonable assumption and sometimes not. Suppose you believe in a multiple regression model as expressed in eqn (6.9),

$$Y = \beta_1 \tau + \beta_2 X_2 + \varepsilon$$

where X_2 is measured without error, but you cannot directly observe τ . Instead you have a measurement of τ (a “fallible indicator”)

$$X_1 = \tau + \delta$$

where δ represents measurement error. And we are particularly interested in the coefficient β_1 (not just whether our multiple regression does a good job of predicting Y).

I will not go through the algebra presented in this section as we will not address the approaches to deal with the subject of measurement error. Suffice to say that, when you compute the regression of Y on X_1 and X_2 , the measurement error in X_1

- “*attenuates*” the coefficient estimate $\hat{\beta}_1$ — makes it systematically biased (smaller in absolute value) with respect to the coefficient β_1 in the true model.
- can bias the coefficient estimate $\hat{\beta}_2$ in either a positive or negative direction (toward or away from 0) depending on sign of covariance between X_1 and X_2 .

Explanation of effect of measurement error in the case of simple linear regression (not from Fox).

$$\begin{aligned} Y_i &= \alpha + \beta_1 \tau_i + \varepsilon_i, \quad X_i = \tau_i + \delta_i \\ &= \alpha + \beta_1 X_i + (\varepsilon_i + \beta_1 \delta_i) = \alpha + \beta_1 \tau_i + \varepsilon_i^* \end{aligned}$$

where $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$, $\delta_i \sim N(0, \sigma_\delta^2)$.

Define $\sigma_\tau^2 = \sum (\tau_i - \bar{\tau})^2 / n$, the variance in the “true” predictor, and assume variability in τ_i is uncorrelated with the measurement error δ_i .

It is an easy exercise to show that

$$E(\hat{\beta}_1) = B_1 / \left(1 + \frac{\sigma_\delta^2}{\sigma_\tau^2} \right) < \beta_1$$

The problem occurs because X_i and ε_i^* are usually not independent for this model

$$\text{Cov}(X_i, \varepsilon_i^*) = \text{Cov}(\tau_i + \delta_i, \varepsilon_i - \beta_1 \delta_i) = -\beta_1 (\sigma_{\tau\delta} + \sigma_\delta^2)$$

Can generally ignore measurement error if

- 1) $\sigma_\delta^2 / \sigma_\tau^2$ is small --- measurement error small relative to true variance in τ 's
- 2) X_i are fixed and predetermined by design ($\tau_i = X_i - \delta_i$, $\sigma_{\tau\delta} + \sigma_\delta^2 = 0$)
- 3) postulated model is $Y_i = \alpha + \beta_1 X_i + \varepsilon_i$

You can probably imagine how measurement error must be a consideration in many social science models where explanatory variables are *concepts* measured with error.

What is “education” in the occupational prestige dataset? Is it a concept measured with error? See Exercise 6.13.

Example of practical importance of measurement error in our current research on effects of air pollution on public health.

- Large EPA-funded study of effects of long-term exposure to “fine particulate matter” ($PM_{2.5}$) and other air pollutants, especially traffic-related pollutants) on cardiovascular disease in older people.
- We want to know the effect of “true exposure” to $PM_{2.5}$ on measures of coronary artery calcium, for example, as well as cardiovascular “events”. This has great implications for public health and regulation of air pollution.
- $PM_{2.5}$ concentration is measured over time at a fixed number of monitoring sites and these data are used to estimate concentrations at the locations where the subjects live. There is considerable error in the estimation of true concentration at the locations where people live (let alone “true exposure”, which involves lots of other factors).
- It is a current research activity to determine good ways to deal with measurement error in this problem.