# 1   Multiple Regression

## 1.1   Motivation: Omitted Variables

Suppose you are interested in understanding how wages are related to years of education, so you look at the model

$$wage = \beta_1 + \beta_2 educ + v.$$

For now, think of $v$ as being the typical error term. The interpretation is that we want to explain *wage* with *educ* and "other stuff" captured in $v$.

Now ask yourself: of the "other stuff" in $v$ that explains wage, is any of that also correlated with education? I am strongly inclined to say "yes." Take IQ for example. I would expect a higher IQ to explain a higher wage; but I also suspect that there is a correlation between IQ and years of education. So when we consider someone with more education, we are also likely considering someone with a higher IQ. This is problematic because $\beta_2$ in the regression above is implicitly telling us the effect of education *and* of IQ on wage, and therefore $\beta_2$ does not isolate the effect of education on wage.

In other words, *we are failing to hold IQ constant when considering different levels of education*, and consequently we are getting both the effect of higher education *and* the effect of higher IQ in our estimate of $\beta_2$. This relationship is illustrated in Figure 1.

That we fail to include a variable that is correlated with both the independent and dependent variable means our estimate for $\beta_2$ will be **biased**, that is, $E[b_2] \neq \beta_2$. We refer to this as **omitted variable bias**. Technically this is consequence of violating classical OLS assumption 2 (see below), i.e. zero conditional mean, because $E[v|educ] \neq 0$.

So how do we progress? Simple: just stick IQ into the regression as well. Our improved model is thus of the form

$$wage = \beta_1 + \beta_2 educ + \beta_3 IQ + u.$$

Now when we take the partial derivative with respect to education, we are explicitly holding IQ constant by definition of a partial derivative. Therefore

$$\frac{\partial wage}{\partial education} = \beta_2$$

gives the relationship between education and wage where IQ is being *controlled for*.

Of course, there are probably other omitted variables as well. In a laboratory experiment, ideally all of these factors can be controlled for if the experiment is properly designed. But
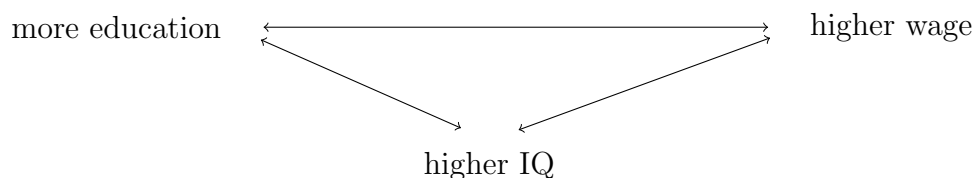
FIGURE 1: More education is correlated with higher wage, but it's also correlated with higher IQ. If we do not hold IQ constant, then we are not accurately characterizing the relationship between education and wage.

we are limited to the data we observe, which may or may not contain all relevant variables. (Probably won't.) Thus, even if we control for a bunch of variables, we still can never be certain that we have fully determined the direct relationship between any $x$ and $y$.

## 1.2   Example: Wages

Import `wages.dta` into Stata. It contains, you guessed it, information about (monthly) wages, education, IQ, and some other stuff. If we regress wages on education, the result is

$$\widehat{wage} = 139.12 + 61.59 \times educ.$$

This implies that someone with one more year of education would be expected to have a higher monthly wage by \$61.586. But as discussed earlier, this is implicitly including the effect of a higher IQ, since the model above fails to control for IQ. We control for IQ by regressing wage on both education and IQ. By doing so, we expect the effect of education to be lower because now the effect isn't being exaggerated by a higher IQ. Indeed,

$$\widehat{wage} = -131.67 + 44.27 educ + 4.95 IQ.$$

So as predicted, the estimated effect of education on wage drops from 61.59 to 44.27. Before controlling for wage, our estimate of $\beta_2$ had an *upward bias*.

The relevant Stata commands and output are found on the following page.

## 2   Classical OLS Assumptions

For OLS to "work" by default, we need the following conditions to hold given dependent variable $y$ and the set of regressors $x_2, x_3, \ldots, x_k$. Note that we have $k-1$ regressors because we started at $x_2$. Therefore we will be estimating $k$ things because we are also estimating the intercept coefficient. Hence we will have $n - k$ degrees of freedom when we start inference.

```
. regress wage educ

      Source |       SS           df       MS        Number of obs   =       852
-------------+----------------------------------      F(1, 850)       =    107.82
       Model | 15622714.1           1 15622714.1      Prob > F        =    0.0000
    Residual |  123165191         850 144900.225      R-squared       =    0.1126
-------------+----------------------------------      Adj R-squared   =    0.1115
       Total |  138787905         851  163088.02      Root MSE        =    380.66


-----------------------------------------------------------------------------
        wage |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+---------------------------------------------------------------
        educ |   61.58627   5.931167    10.38   0.000     49.94482    73.22772
       _cons |   139.1171   81.16017     1.71   0.087    -20.18069     298.415
-----------------------------------------------------------------------------

. regress wage educ iq

      Source |       SS           df       MS        Number of obs   =       852
-------------+----------------------------------      F(2, 849)       =     67.17
       Model | 18960227.2           2 9480113.62      Prob > F        =    0.0000
    Residual |  119827678         849 141139.786      R-squared       =    0.1366
-------------+----------------------------------      Adj R-squared   =    0.1346
       Total |  138787905         851  163088.02      Root MSE        =    375.69


-----------------------------------------------------------------------------
        wage |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+---------------------------------------------------------------
        educ |   44.26802   6.851945     6.46   0.000     30.81928    57.71676
          iq |   4.954005   1.018755     4.86   0.000     2.954432    6.953578
       _cons |  -131.6712    97.5547    -1.35   0.177    -323.1479    59.80543
-----------------------------------------------------------------------------
```

1. **MLR1: Correct Linear Model.** The true model is linear and correctly specified as

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \beta_k x_k + u. \tag{1}$$

Intuition: if we estimate a population model that's actually of a different form, then our estimates are probably garbage.

2. **MLR2: Zero Conditional Mean.** The error term has zero mean conditional upon the regressors, that is,

$$E[u|x_2, \ldots, x_k] = 0. \tag{2}$$

Intuition: think of the error term as being the mistake of the model. If we expect the mistake to be non-zero on average, then our model is probably garbage. This condition is equivalent to saying that $u$ is uncorrelated with all of the regressors.

More technically, it allows us to go from

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \beta_k x_k + u$$

to the conditional expectation

$$E[y|x_2, \ldots, x_k] = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \beta_k x_k,$$

the latter being our equation for fitted values and predictions of $y$.

3. **MLR3: Homoskedasticity.** The conditional variance of the error term is constant and finite, that is,

$$\mathrm{Var}(u|x_2, \ldots, x_k) = \sigma_u^2 < \infty. \tag{3}$$

There isn't much economic intuition here; it's mostly a technical assumption, albeit an unrealistic one, that offers a convenient starting point for rigorous analysis. In practice it is violated frequently, which is not difficult to deal with (as explained later). This condition is illustrated in Figure 2.

4. **MLR4: Independent Errors.** Errors for different observations are statistically independent, that is,

$$u_i \perp u_j \quad \text{whenever } i \neq j.$$

Intuition: if model errors are correlated, then there is some underlying pattern that we are overlooking, so our results are probably garbage.
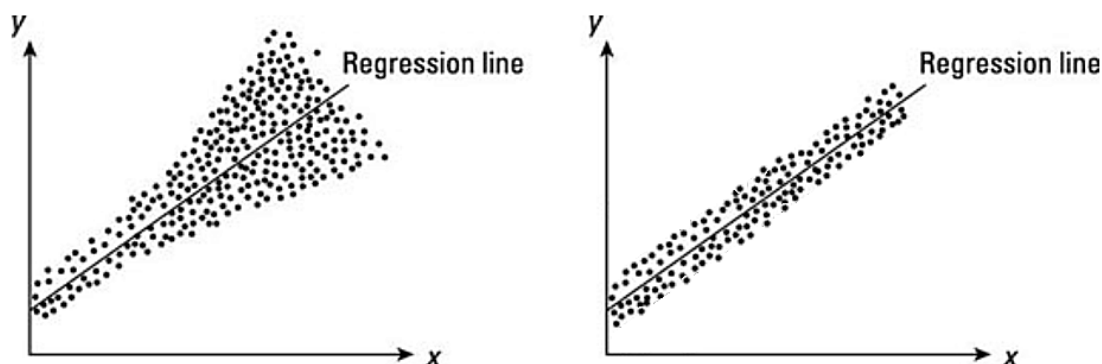


FIGURE 2: The figure on the left is an example of heteroskedasticity; the right an example of homoskedasticity. The left is heteroskedastic because the variation around the regression line gets bigger as $x$ increases. Good luck envisioning this in higher dimensions.

As an example of a violation, suppose we look at ECN 102 final exam scores in all of 2017; that means we're looking at ECN 102 final exam scores for three different professors. Problem is, different professors write exams of differing difficulty. Hence we would expect a lenient professor's students to do better than the regression predicts (so we'd have correlation among observations with positive $u$), and we expect a challenging professor's students to do worse than what the regression predicts (so we'd have correlation among observations with negative $u$). This is called **clustering** because each professor's final exam forms a cluster of students.

5. **MLR5: Normality of Errors.** Errors are normally distributed with some variance $\sigma^2$, that is,

$$u_i \sim \mathcal{N}\left(0, \sigma^2\right). \tag{4}$$

This is another technical assumption for "nice" results, explained below. In practice it can be weakened.

6. **MLR6: No Perfect Multicollinearity.** There exists no exact linear relationship between explanatory samples. Furthermore, the number of observations must be greater than the number of explanatory variables (plus constant term), i.e. $n \geq k$.

Intuition: if there is such a perfect relationship between two or more regressors, then we can't "untangle" the effect of each regressor. In other words, it's like including the same regressor twice, and that redundancy breaks the OLS solution technique.

Assumptions MLR1-2 imply that OLS estimates are unbiased, so that $E[b_j] = \beta_j$. Assumptions MLR1-4 imply that OLS estimates are consistent, so that $b_j \xrightarrow{p} \beta_j$ as $n \to \infty$. Furthermore, assumptions MLR1-4 imply that OLS is the **b**est **l**inear **u**nbiased **e**stimator, or BLUE. When we say "best," we mean we have the smallest standard errors and hence precision of inference is the most accurate. If we throw in assumption MLR5, then OLS is the **b**est **u**nbiased **e**stimator, even when compared to nonlinear methods. (Note that assumption MLR5 is needed to do inference with small sample sizes.) Assumption MLR6 is absolutely required; in the presence of perfect multicollinearity, the regression cannot be executed. Accordingly, this is usually just implicitly assumed because otherwise it's game over and we should just give up and go home.

# 3   Multiple Regression Inference

Under MLR1-4, the $t$-statistic regarding regressor $x_j$ is given by

$$t = \frac{b_j - \beta_j}{se(b_j)},$$ (5)

and it is drawn from an approximate $T(n - k)$ distribution. (We are estimating $k$ things, which is why we have $n - k$ degrees of freedom.) Inference proceeds in the usual way.

There is no rule of thumb for how large $n$ needs to be for the approximation to be adequate. If MLR5 holds, then $t$ is drawn from exact $T(n - k)$ distribution. If either MLR3 or MLR4 fail, then the typical standard errors are not valid. Instead we need to use, respectively, **heteroskedasticity-robust standard errors** or **cluster-robust standard errors**. These are both very easy to implement in Stata, as shown in a forthcoming example.

# 4   Dummy Variables

## 4.1   Definition of Dummy Variable

We might be interested in seeing how different categories affect the dependent variable. For instance, we might want to see if someone working in an urban environment earns more than someone working elsewhere. To analyze, we construct a **dummy variable** that is equal to either zero or one. An urban worker would have value $urban = 1$, and a non-urban worker would have value $urban = 0$. Accordingly, we would run the regression

$$wage = \beta_1 + \beta_2 educ + \beta_3 IQ + \beta_4 urban + u.$$

The coefficient $\beta_4$, then, would tell you the expected difference in monthly wage for an urban worker compared to a non-urban worker. Another way of thinking about it is, $\beta_4$ captures the expected change in wage if a worker moves from a non-urban environment to an urban environment, that is, if $urban$ changes from 0 to 1.

## 4.2   Dummy Variable Trap

Notice in the preceding example that there are two categories, but only one dummy variable. In general, if you have $m$ categories, then you must include exactly $m - 1$ dummy variables; the category you omit is called the **reference category**. Including dummy variables for all possible categories results in the **dummy variable trap**, which is a source of perfect

multicollinearity that breaks OLS estimation. So always use one fewer dummy than there are categories.

Here's a really stupid example to illustrate why things go wrong. Suppose everyone has a choice of either having Swedish Fish, Sour Patch Kids, or Mike and Ikes, but can only choose one. We want to see how many cavities each person receives from eating so much damn candy. We record their choices in the following manner:

$$choice = 1 \quad \text{if Swedish Fish,}$$

$$choice = 2 \quad \text{if Sour Patch Kids,}$$

$$choice = 3 \quad \text{if Mike and Ikes.}$$

Now let's create dummies for all categories. Let $d_1 = 1$ for choosing Swedish Fish; $d_2 = 1$ for choosing Sour Patch Kids; and $d_3 = 1$ for choosing Mike and Ikes. Then the possible values for each dummy are

$$choice = 1 \quad \implies \quad d_1 = 1, \; d_2 = 0, \; d_3 = 0,$$

$$choice = 2 \quad \implies \quad d_1 = 0, \; d_2 = 1, \; d_3 = 0,$$

$$choice = 3 \quad \implies \quad d_1 = 0, \; d_2 = 0, \; d_3 = 1.$$

Notice that in all three cases, $d_1 + d_2 + d_3 = 1$. And therefore, say, $d_1 = 1 - d_2 - d_3$. This is perfect multicollinearity because one of our regressors $(d_1)$ can be perfectly explained by a linear relationship of two other regressors $(d_2$ and $d_3)$. So if we try to regress *cavities* on $d_1$, $d_2$, and $d_3$, then OLS explodes and we're all doomed.

Except you can just remove any one of the three dummies from the regression, then all is well and well is all for all. The coefficients of the model are then seen as being *relative* to the reference category. To that end, consider the model where we omit the Swedish Fish dummy variable $d_1$, given by

$$cavities = \beta_1 + \beta_2 d_2 + \beta_3 d_3 + u.$$

Let us interpret each coefficient.

- $\beta_1$: how many cavities are associated with eating Swedish Fish (the reference category);

- $\beta_2$: how many more (or less, if negative) cavities are associated with eating Sour Patch Kids instead of Swedish Fish;

- $\beta_3$: how many more (or less, if negative) cavities are associated with eating Mike and Ikes instead of Swedish Fish.

In the case of the urban workers, $\beta_4$ captures how much higher of a wage a person receives if they work in an urban environment relative to working in a non-urban environment (the reference category).

## 4.3   Example: Wages

Again using `wages.dta`, let us consider the regression proposed earlier,

$$wage = \beta_1 + \beta_2 educ + \beta_3 IQ + \beta_4 urban + u.$$

OLS estimation yields

$$\widehat{wage} = -213.28 + 41.58 educ + 4.92 IQ + 169.01 urban.$$

The $p$-value for $\beta_4$ indicates statistically significance, so we conclude that an urban worker is expected to earn a monthly wage \$169.01 higher than that of a non-urban worker. To account for the possibility of heteroskedasticity, I tell Stata to use heteroskedasticity-robust standard errors with the option `vce(robust)`.

```
. regress wage educ iq urban, vce(robust)

Linear regression                               Number of obs   =        852
                                                F(3, 848)       =      53.41
                                                Prob > F        =     0.0000
                                                R-squared       =     0.1719
                                                Root MSE        =     368.15

-------------------------------------------------------------------------------
             |               Robust
        wage |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+-----------------------------------------------------------------
        educ |   41.58144   6.793912     6.12   0.000     28.24658    54.91629
          iq |   4.919558    .944874     5.21   0.000     3.064992    6.774124
       urban |   169.0137   26.54763     6.37   0.000      116.907    221.1205
       _cons |  -213.2816   95.91454    -2.22   0.026    -401.5393   -25.02381
-------------------------------------------------------------------------------
```