

POLS 503: Advanced Quantitative Political
Methodology: The Notes

Jeffrey B. Arnold

2016-04-19

Contents

1	Introduction	5
2	Linear Regression and the Ordinary Least Squares (OLS) Estimator	7
2.1	Linear Regression Function	7
2.2	Ordinary Least Squares	8
2.3	Properties of the OLS Estimator	9
2.4	References	10
3	OLS Troubleshooting and Diagnostics	11
3.1	Multi-Collinearity	11
3.2	Omitted Variable Bias	11
3.3	Measurement Error	11
3.4	Non-linearity	11
3.5	Heteroskedasticity and Auto-correlation	11
3.6	Non-constant variance (Heteroskedasticity)	12
4	Appendix	13
4.1	Multivariate Normal Distribution	13

Chapter 1

Introduction

hello, world!

Chapter 2

Linear Regression and the Ordinary Least Squares (OLS) Estimator

Since we will largely be concerned with using linear regression for inference, we will start by discussion the population parameter of interest (population linear regression function), then the sample statistic (sample linear regression function) and estimator (ordinary least squares).

We will then consider the properties of the OLS estimator.

2.1 Linear Regression Function

The **population linear regression function** is

$$r(x) = E[Y|X = x] = \beta_0 + \sum_{k=1}^K \beta_k x_k.$$

The population linear regression function is defined for random variables, and will be the object to be estimated.

Names for \mathbf{y}

- dependent variable
- explained variable
- response variable
- predicted variable
- regressand
- outcome variable

Names for \mathbf{X} ,

- independent variables
- explanatory variables
- treatment and control variables
- predictor variables
- covariates
- regressors

To estimate the unknown population linear regression, we will use the **sample linear regression function**,

$$\hat{r}(x_i) = \hat{y}_i = \hat{\beta}_0 + \sum_{k=1}^K \hat{\beta}_k x_k.$$

However, we

\hat{Y}_i are the fitted or predicted value The **residuals** or **errors** are the prediction errors of the estimates

$$\hat{\epsilon}_i = y_i - \hat{y}_i$$

β are the parameters; β_0 is called the *intercept*, and β_1, \dots, β_K are called the *slope parameters*, or *coefficients*.

We will then consider the properties of the OLS estimator.

The linear regression can be more compactly written in matrix form,

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} & \cdots & x_{K,1} \\ 1 & x_{1,2} & x_{2,2} & \cdots & x_{K,2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1,N} & x_{2,n} & \cdots & x_{K,N} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}.$$

More compactly, the linear regression model can be written as,

$$\underbrace{\mathbf{y}}_{N \times 1} = \underbrace{\mathbf{X}}_{N \times K} \underbrace{\boldsymbol{\beta}}_{K \times 1} + \underbrace{\boldsymbol{\epsilon}}_{N \times 1}.$$

The matrix \mathbf{X} is called the *design* matrix. Its rows are each observation in the data. Its columns are the intercept, a column vector of 1's, and the values of each predictor.

2.2 Ordinary Least Squares

Ordinary least squares (OLS) is an estimator of the slope and statistic of the regression line¹. OLS finds values of the intercept and slope coefficients by minimizing the squared errors,

$$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_K = \arg \min_{b_0, b_1, \dots, b_K} \sum_{i=1}^N \underbrace{\left(y_i - b_0 - \sum_{k=1}^K b_k x_{i,k} \right)^2}_{\text{squared error}},$$

or, in matrix notation,

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \arg \min_{\mathbf{b}} \sum_{i=1}^N (y_i - \mathbf{b}' \mathbf{x}_i)^2 \\ &= \arg \min_{\mathbf{b}} \sum_{i=1}^N u_i^2 \\ &= \arg \min_{\mathbf{b}} \mathbf{u}' \mathbf{u} \end{aligned}$$

where $\mathbf{u} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$.

In most statistical models, including even generalized linear models such as logit, the solution to this minimization problem would be solved with optimization methods that require iteration. One nice feature

¹Ordinary least squares is distinguished from *generalized least squares* (GLS).

of OLS is that there is a closed form solution for $\hat{\beta}$ even in the multiple regression case, so no iterative optimization methods need to be used.

In the bivariate regression case, the OLS estimators for β_0 and β_1 are

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^N (x_i - \bar{x})^2} \\ &= \frac{\text{Cov}(\mathbf{x}, \mathbf{y})}{\text{Var } \mathbf{x}} = \frac{\text{Sample covariance between } \mathbf{x} \text{ and } \mathbf{y}}{\text{Sample variance of } \mathbf{x}}.\end{aligned}$$

In the multiple regression case, the OLS estimator for $\hat{\beta}$ is

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}.$$

The term $\mathbf{X}'\mathbf{X}$ is similar to the variance of \mathbf{x} in the bivariate case. The term $\mathbf{X}'\mathbf{y}$ is similar to the covariance between \mathbf{X} and \mathbf{y} in the bivariate case.

The sample linear regression function estimated by OLS has the following properties:

1. Residuals sum to zero,

$$\sum_{i=1}^N \hat{\epsilon}_i = 0.$$

This implies that the mean of residuals is also 0.

2. The regression function passes through the point $(\bar{\mathbf{y}}, \bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_K)$. In other words, the following is always true,

$$\bar{\mathbf{y}} = \hat{\beta}_0 + \sum_{k=1}^K \hat{\beta}_k \bar{\mathbf{x}}_k.$$

3. The residuals are uncorrelated with the predictor

$$\sum_{i=1}^N x_i \hat{\epsilon}_i = 0$$

4. The residuals are uncorrelated with the fitted values

$$\sum_{i=1}^N \hat{y}_i \hat{\epsilon}_i = 0$$

2.3 Properties of the OLS Estimator

2.3.1 What makes an estimator good?

Estimators are evaluated not on how close an estimate in a given sample is to the population, but how their sampling distributions compare to the population. In other words, judge the *methodology* (estimator), not the *result* (estimate).^[^ols-properties-references]

Let θ be the population parameter, and $\hat{\theta}$ be an estimator of that population parameter.

Bias The bias of an estimator is the difference between the mean of its sampling distribution and the population parameter,

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta.$$

Variance The variance of the estimator is the variance of its sampling distribution, $\text{Var}(\hat{\theta})$.

Efficiency (Mean squared error) An efficient estimator is one that minimizes a given “loss function”, which is a penalty for missing the population average. The most common loss function is squared loss, which gives the *Mean Squared Error (MSE)* of an estimator.

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = (E(\hat{\theta}) - \theta)^2 + E(\hat{\theta} - E(\hat{\theta}))^2 = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$$

The mean squared error is a function of both the bias and variance of an estimator.

This means that some biased estimators can be more efficient : than unbiased estimators if their variance offsets their bias.²

Consistency is an asymptotic property³, that roughly states that an estimator converges to the truth as the number of observations grows, $E(\hat{\theta} - \theta) \rightarrow 0$ as $N \rightarrow \infty$. Roughly, this means that if you had enough (infinite) data, the estimator will give you the true value of the parameter.

2.3.2 Properties of OLS

- When is OLS unbiased?
- When is OLS consistent?
- When is OLS efficient?

Assumption	Formal statement	Consequence of violation
No (perfect) collinearity	$\text{rank}(\mathbf{X}) = K, K < N$	Coefficients unidentified
\mathbf{X} is exogenous	$E(\mathbf{X}\boldsymbol{\varepsilon}) = 0$	Biased, even as $N \rightarrow \infty$
Disturbances have mean 0	$E(\boldsymbol{\varepsilon}) = 0$	Biased, even as $N \rightarrow \infty$
No serial correlation	$E(\varepsilon_i \varepsilon_j) = 0, i \neq j$	Unbiased, wrong se
Homoskedastic errors	$E(\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon})$	Unbiased, wrong se
Gaussian errors	$\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2)$	Unbiased, se wrong unless $N \rightarrow \infty$

Note that these assumptions can be sometimes be written in largely equivalent, but slightly different forms.

2.4 References

- Wooldrige, Ch 3.
- Fox, Ch 6, 9.

²It follows from the definition of MSE, that biased estimator, $\hat{\theta}_B$, has a lower MSE than an unbiased estimator, $\hat{\theta}_U$, if $\text{Bias}(\hat{\theta}_B)^2 < \text{Var}(\hat{\theta}_U) - \text{Var}(\hat{\theta}_B)$.

³As the number of observations goes to infinity.

Chapter 3

OLS Troubleshooting and Diagnostics

3.1 Multi-Collinearity

3.2 Omitted Variable Bias

3.3 Measurement Error

3.4 Non-linearity

3.5 Heteroskedasticity and Auto-correlation

Note, that OLS assumes that the variance of the the disturbances is constant $\hat{Y} - Y = \varepsilon = \sigma^2$. What happens if it isn't?

$$\mathbf{\Sigma} = \begin{bmatrix} \text{Var}(\varepsilon_1) & \text{Cov}(\varepsilon_1, \varepsilon_2) & \cdots & \text{Cov}(\varepsilon_1, \varepsilon_N) \\ \text{Var}(\varepsilon_2, \varepsilon_1) & \text{Var}(\varepsilon_2) & \cdots & \text{Cov}(\varepsilon_2, \varepsilon_N) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\varepsilon_N, \varepsilon_1) & \text{Cov}(\varepsilon_N, \varepsilon_2) & \cdots & \text{Cov}(\varepsilon_N) \end{bmatrix} \Sigma = \begin{bmatrix} \text{E}(\varepsilon_1^2) & \text{E}(\varepsilon_1 \varepsilon_2) & \cdots & \text{E}(\varepsilon_1 \varepsilon_N) \\ \text{E}(\varepsilon_2 \varepsilon_1) & \text{E}(\varepsilon_2^2) & \cdots & \text{E}(\varepsilon_2 \varepsilon_N) \\ \vdots & \vdots & \ddots & \vdots \\ \text{E}(\varepsilon_N \varepsilon_1) & \text{E}(\varepsilon_N \varepsilon_2) & \cdots & \text{E}(\varepsilon_N^2) \end{bmatrix}$$

The matrix can be written more compactly as,

$$\mathbf{\Sigma} = \text{E}(\boldsymbol{\varepsilon} \boldsymbol{\varepsilon}')$$

An assumption is that errors are independent, $\text{E}(\varepsilon_i \varepsilon_j)$ for all $i \neq j$. This means that all off-diagonal elements of $\mathbf{\Sigma}$ are 0\$. Additionally, all ε_i are assumed to have the same variance, σ^2 . Thus, the variance-covariance matrix of the errors is assumed to have a diagonal matrix with the form,

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}_N$$

If these assumptions of the errors do not hold, then Σ does not take this form, and more complicated models than OLS need to be used to get correct standard errors.

3.6 Non-constant variance (Heteroskedasticity)

The homoskedastic case assumes that each error term has its own variance. In the heteroskedastic case, each disturbance may have its own variance, but they are still uncorrelated (Σ is diagonal)

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N^2 \end{bmatrix}$$

The problem is that now there are N variance parameters to estimate, in addition to the K slope coefficients. Now, there are more parameters than we can estimate. With heteroskedasticity, OLS will be unbiased, but the standard errors will be incorrect.

More general case allows for heteroskedasticity, and autocorrelation ($\text{Cov}(\varepsilon_i, \varepsilon_j) \neq 0$),

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} & \cdots & \sigma_{1,N} \\ \sigma_{2,1} & \sigma_2^2 & \cdots & \sigma_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{N,1} & \sigma_{N,2} & \cdots & \sigma_N^2 \end{bmatrix}$$

As with heteroskedasticity, OLS will be unbiased, but the standard errors will be incorrect.

Chapter 4

Appendix

4.1 Multivariate Normal Distribution

The multivariate normal distribution is the generalization of the univariate normal distribution to more than one dimension.¹ The random variable, \mathbf{x} , is a length k vector. The k length vector $\boldsymbol{\mu}$ are the means of \mathbf{x} , and the $k \times k$ matrix, $\boldsymbol{\Sigma}$, is the variance-covariance matrix,

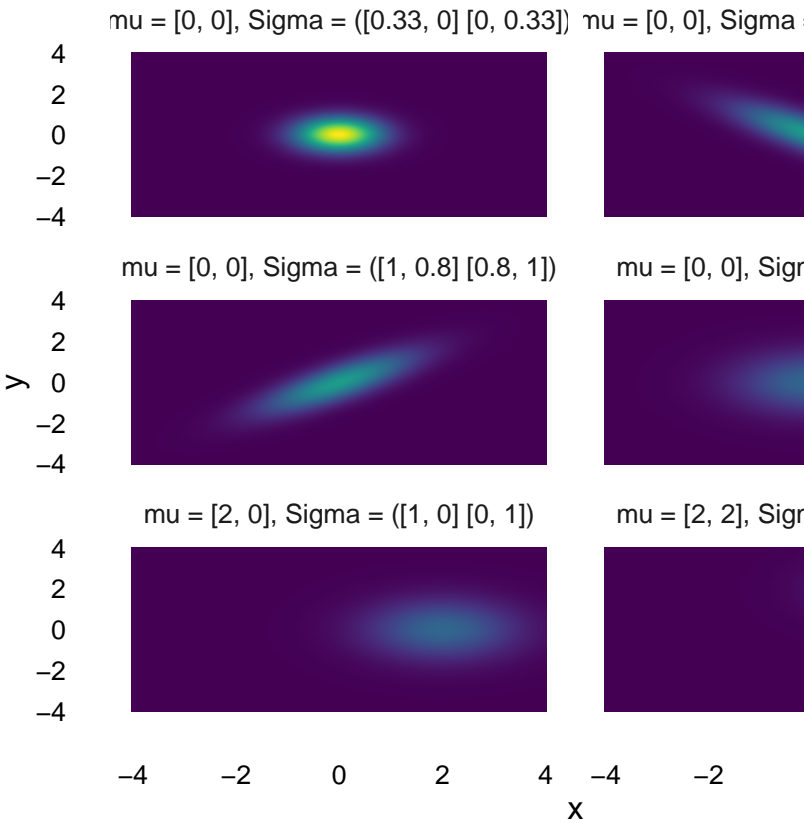
$$\mathbf{x} \sim \mathcal{N}_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \end{bmatrix} \sim \mathcal{N}_k \left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \sigma_{1,2} & \cdots & \sigma_{1,k} \\ \sigma_{2,1} & \sigma_2^2 & \cdots & \sigma_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{k,1} & \sigma_{k,2} & \cdots & \sigma_{k,k} \end{bmatrix} \right)$$

The density function of the multivariate normal is,

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2k)^{-\frac{k}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right).$$

You can sample from and calculate the density for the multivariate normal distribution with the functions `dmvnorm` and `rmvnorm` from the package **mvtnorm**.

¹See Multivariate normal distribution and references therein.



Density plots of different bivariate normal distributions,