

Chapter 7

Linear and generalised linear models

7.1 Linear model

7.1.1 Model

So far y_i has been *iid*. Here we extend this to the case where the observations depend linearly on some known (fixed) $p \times 1$ vector of regressors x_i . The basic structure will be that

$$y_i = x_i' \beta + \varepsilon_i, \quad E(\varepsilon_i) = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2,$$

and the ε_i are iid. This model is indexed by the parameters β and σ^2 . We can write the model in matrix form as $y = X\beta + \varepsilon$, where $y = (y_1, \dots, y_n)'$, $X = (x_1, \dots, x_n)'$ and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)'$.

7.2 Definition

$$\hat{\beta} = (X'X)^{-1} X'y, \quad \text{assuming } X \text{ has full rank.}$$

7.2.1 Interpretation

- (1) The most common association with the above estimator is with the so called least squares principle. This write

$$\begin{aligned} \hat{\beta} &= \arg \min_{\beta} \varepsilon' \varepsilon = \arg \min_{\beta} \sum_{i=1}^n \varepsilon_i^2 \\ &= \arg \min_{\beta} \sum_{i=1}^n (y_i - x_i' \beta)^2. \end{aligned}$$

That is it determines $\hat{\beta}$ by finding the value of β which minimises the squared vertical distances between y_i and $x_i' \beta$. This construction leaves many questions unanswered for the choice of squared distances looks completely ad hoc. Justifications usually involve the resulting estimator has good properties. Alternatives include

$$\tilde{\beta} = \arg \min_{\beta} \sum_{i=1}^n |y_i - x_i' \beta|,$$

which is called an L_1 (or least absolute deviation) regression estimator, or

$$\tilde{\tilde{\beta}} = \arg \min_{\beta} \text{medium} \{ |y_1 - x_1' \beta|, \dots, |y_n - x_n' \beta| \},$$

the least medium absolute deviation (high breakdown) estimator.

- (2) Suppose we regard X as fixed (or we condition on it for purposes of inference) and $y_i \sim N(x_i' \beta, \sigma^2)$, with the y_i being independent over i . Then the ML estimator is

$$\begin{aligned} \hat{\beta} &= \arg \max_{\beta} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i' \beta)^2 \right\} \\ &= \arg \min_{\beta} \sum_{i=1}^n (y_i - x_i' \beta)^2, \end{aligned}$$

the same as the least squares estimator. This implies least squares estimation is the ML estimator if the model is Gaussian and has many optimality properties if the truth is a Gaussian regression.

- (3) Suppose $y_i \sim \text{Laplace}(x_i' \beta, \sigma^2)$, that is

$$f(y_i; \beta, \sigma^2) = \frac{1}{2\sigma} \exp \left\{ -\frac{1}{\sigma} |y_i - x_i' \beta| \right\},$$

and that the data are independent over i . For this distribution the usual mean and variance of the Gaussian regression carry over. Then the model has much fatter tails in $y_i - x_i' \beta$ than the normal as the Laplace regression model involves only exponentiating minus the absolute value, while the Gaussian model squares these distances. The ML estimator for β is

$$\begin{aligned} \tilde{\beta} &= \arg \max_{\beta} \left(\frac{1}{2\sigma} \right)^n \exp \left\{ -\frac{1}{\sigma} \sum_{i=1}^n |y_i - x_i' \beta| \right\} \\ &= \arg \min_{\beta} \sum_{i=1}^n |y_i - x_i' \beta|. \end{aligned}$$

That is the ML estimator for this distribution is an L_1 regression. Such a regression is typically less sensitive to unusual observations than the least squares regression — this can be a good thing or a bad thing in practice. NOTE: the Laplace log-likelihood function is not continuously differentiable with respect to β and so many of the typical properties of MLE do not go through even if the data is Laplace.

example

Think about the dataset $x = (1, 2, 3, 5)'$ and $y = (1, 2, 3, 1)'$, so 3 observations are in a straight-line with gradient one and a single observation is an outlier to that. The LS estimator is $(1 + 4 + 9 + 5) / (1 + 4 + 9 + 25) = 0.487$ and is deeply effected by the outlier. The L_1 estimator minimises sum of the absolute values. These decrease with a gradient of -4 until $\beta = 1/5$ when it slows to -2 which continues until $\beta = 1$ when the gradient switches to 4. So the L_1 regressor gives an estimator of 1, so ignores the outlier completely.

7.3 Properties of least squares

For $y = X\beta + \varepsilon$ and $E(\varepsilon) = 0$ while $\text{Var}(\varepsilon) = \sigma^2 I$, then

$$\hat{\beta} = (X'X)^{-1} X'y = \beta + (X'X)^{-1} X'\varepsilon,$$

so $\hat{\beta}$ is unbiased while $\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$.

If ε is Gaussian, then $\hat{\beta}$ is a linear combination of Gaussian and so if Gaussian. Further,

$$\frac{\partial \log L}{\partial \beta} = \frac{1}{\sigma^2} X' (y - X\beta) \quad \text{implying} \quad \text{Var} \frac{\partial \log L}{\partial \beta} = \frac{1}{\sigma^2} X' X,$$

so $\hat{\beta}$ achieves the Cramer-Rao lower bound and is the minimum variance unbiased estimator.

If the parametric model is not Gaussian $\hat{\beta}$ has more limited optimality properties. In particular it is possible to show that in the class of estimators which is linear in y , $\hat{\beta}$ is the minimum variance unbiased estimator (Gauss-Markov theorem).

This result is less helpful than it at first seems for when the data is non-Gaussian it seems clear that we should be using a non-linear estimator! See, for example, the discussion above about the Laplace regression model.

A more thorough discussion of these issues is given in Chapter ??.

7.4 Generalized linear models*

7.4.1 Models

Instead of working with a particular density function it is sometimes helpful to construct a class of densities. One such class is the exponential family, which puts

$$f(y; \theta, \psi) = \exp \left[\frac{1}{\psi} \{y\theta - b(\theta)\} - c(y, \psi) \right].$$

We are free to chose the functions $b(\cdot)$ and $c(\cdot, \cdot)$ to place familiar distributions in this framework.. Simple examples of this includes

- $N(\mu, \sigma^2)$, $\theta = \mu$, $\psi = \sigma^2$, $b(\theta) = \theta^2/2$.
- Poisson with mean μ , puts $\theta = \log \mu$, $b(\theta) = \exp \theta$ and $\psi = 1$.
- Bernoulli with probability of success of p . Then we put $\theta = \log \{p/(1-p)\}$, $b(\theta) = \log(1 + \exp \theta)$ and $\psi = 1$.

Other common distributions which go into this framework includes the gamma and some extreme value distributions. An excellent review of this material is given in Azzalini (1996, Ch. 6) and Garthwaite, Jolliffe and Jones (1995, Ch. 10).

Having a class of densities is useful as it allows us to derive some generic results. First, the log-likelihood is

$$\log L(\theta, \psi; y) = \text{const} + \frac{\theta}{\psi} \sum_{i=1}^n y_i - \frac{n b(\theta)}{\psi} - \sum c(y_i, \psi),$$

which implies

$$\frac{\partial \log L(\theta, \psi; y)}{\partial \theta} = \frac{1}{\psi} \sum_{i=1}^n y_i - \frac{n b'(\theta)}{\psi}, \quad \text{where} \quad b'(\theta) = \frac{\partial b(\theta)}{\partial \theta}.$$

As the expectation of the score is zero when evaluated at the true parameter point this implies

$$E(y_i) = b'(\theta),$$

while the information equality implies

$$\text{Var}(y_i) = \psi b''(\theta).$$

7.4.2 Regression models

The Gaussian distribution $Y_i \sim N(\mu, \sigma^2)$ extends to the regression problem where $Y_i \sim N(\beta'x_i, \sigma^2)$. An interesting question is how this argument extends to non-Gaussian densities. Clearly we cannot just write $\mu_i = E(y_i) = x_i'\beta$ as some non-Gaussian densities require the mean to be positive! Instead we will work with a function of the mean

$$g(\mu_i) = x_i'\beta = \eta_i,$$

which is called a link function. Throughout I will assume $g(\cdot)$ is a monotonic function. An example of this is the Poisson case where $\mu_i > 0$ and so it makes some sense to write $\log \mu_i = x_i'\beta$. Economic examples of the use of Poisson type regressions includes models of patents by Hausman, Hall and Griliches (1984). Likewise in the Bernoulli case we have $\mu_i \in (0, 1)$ and has been used by, for example, Micklewright (1989) in labour economics. Hence it might make sense to use a logistic transformation $\log \{\mu_i / (1 - \mu_i)\} = x_i'\beta$.

If we combine an exponential family model with a regression model then we call the result a generalized linear model. It has the generic form

$$\log f(y; \theta_i, \psi) = \frac{1}{\psi} \{y\theta_i - b(\theta_i)\} - c(y, \psi),$$

where

$$\mu_i = E(y_i) = b'(\theta_i) \quad \text{and} \quad g(\mu_i) = x_i'\beta = \eta_i.$$

Hence the likelihood depends now only on the parameters ψ and β as the θ_i are solely determined by the regressors and the regression coefficients.

7.4.3 Likelihood analysis

We think of β as the parameters of interest, then

$$\log L(\beta, \psi; y) = \frac{1}{\psi} \sum_{i=1}^n \{y_i \theta_i - b(\theta_i)\} - \sum_{i=1}^n c(y_i, \psi),$$

and so

$$\frac{\partial \log L(\beta, \psi; y)}{\partial \beta} = \frac{1}{\psi} \sum_{i=1}^n \left[\frac{\partial \eta_i}{\partial \beta} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \theta_i}{\partial \mu_i} \{y_i - b'(\theta_i)\} \right].$$

Although this looks complicated it simplifies to the important result

$$\frac{\partial \log L(\beta, \psi; y)}{\partial \beta} = \sum_{i=1}^n \left[x_i \frac{\{y_i - \mu_i\}}{\text{Var}(y_i)} \frac{\partial \mu_i}{\partial \eta_i} \right],$$

as

$$\frac{\partial \eta_i}{\partial \beta} = x_i, \quad \frac{\partial \mu_i}{\partial \theta_i} = b''(\theta_i) = \frac{\text{Var}(y_i)}{\psi}.$$

This score cannot be immediately solved to deliver the ML estimator of β as in general $\frac{\partial \mu_i}{\partial \eta_i}$, $\text{Var}(y_i)$ and μ_i all depend on β . Hence we will have to use a numerical maximisation procedure to iterate to the maximum. Here we will look at a Fisher scoring method¹, which uses the expected information matrix.

The expected information for the sample can be found by looking at

$$E \left\{ \frac{\partial \log L(\beta, \psi; Y_i)}{\partial \beta_j} \frac{\partial \log L(\beta, \psi; Y_i)}{\partial \beta_k} \right\} = E \left[x_{ij} \frac{\{y_i - \mu_i\}}{\text{Var}(y_i)} \frac{\partial \mu_i}{\partial \eta_i} x_{ik} \frac{\{y_i - \mu_i\}}{\text{Var}(y_i)} \frac{\partial \mu_i}{\partial \eta_i} \right],$$

which simplifies to

$$x_{ij} x_{ik} \frac{1}{\text{Var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2.$$

This implies, using independence across the Y_i that

$$-E \left(\frac{\partial^2 \log L}{\partial \beta \partial \beta'} \right) = E \left\{ \frac{\partial \log L}{\partial \beta} \frac{\partial \log L}{\partial \beta'} \right\} = \sum_{i=1}^n x_i x_i' \frac{1}{\text{Var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2 = X' W X,$$

where W is diagonal with $ii - th$ element

$$w_i = \frac{1}{\text{Var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2,$$

and $X_{ij} = x_{ij}$. Likewise the score can be written as

$$\frac{\partial \log L(\beta, \psi; y)}{\partial \beta} = X' W u, \quad \text{where} \quad u_i = \{y_i - \mu_i\} \frac{\partial \mu_i}{\partial \eta_i}.$$

The implication of this is that the Fisher scoring algorithm actually carries out iteratively weighted least squares, computing the update of β by

$$\beta^{(k+1)} = \beta^{(k)} + (X' W X)^{-1} X' W u = (X' W X)^{-1} X' W (X \beta^{(k)} + u) = (X' W X)^{-1} X' W y^{(k)},$$

say.

¹Fisher scoring generically iterates

$$\theta^{(k+1)} = \theta^{(k)} + \left(-E \frac{\partial^2 \log L}{\partial \theta \partial \theta'} \Big|_{\theta=\theta^{(k)}} \right)^{-1} \frac{\partial \log L(\theta^{(k)}, y)}{\partial \theta}$$

Hence it is a particular case of the quasi-Newton method.

exercise

Suppose we wish to model durations of employment, y_i , as a function of some (assumed fixed) explanatory variable, x_i . We suppose

$$f(y_i; \beta) = \lambda_i \exp(-y_i \lambda_i), \quad y_i > 0, \quad \lambda_i = \exp(x_i \beta),$$

and $E(y_i) = \lambda_i^{-1}$. Assume the y_i are independent over i . Write down the model's likelihood and find the score function for β . Show that the log-likelihood is concave and hence suggest how you might numerically compute the ML estimator?

7.4.4 Asymptotic distribution

In our discussion of the asymptotic distribution of the ML estimator we have always assumed that the observations are *iid*. This model does not have identical distributions, although the independence across individuals remains.

The average expected information per observation is

$$\frac{1}{n} E \left(\frac{\partial \log L}{\partial \beta} \frac{\partial \log L}{\partial \beta'} \right) = \frac{1}{n} \sum_{i=1}^n E \left(\frac{\partial \log L(\theta; Y_i)}{\partial \beta} \frac{\partial \log L(\theta; Y_i)}{\partial \beta'} \right) = \frac{1}{n} X' W X.$$

Now suppose that

$$\lim_{n \rightarrow \infty} \frac{1}{n} E \left(\frac{\partial \log L}{\partial \beta} \frac{\partial \log L}{\partial \beta'} \right) = M > 0,$$

then we have that

$$\sqrt{n} (\hat{\beta} - \beta) \xrightarrow{d} N(0, M^{-1}).$$

This result will hold as long as M is finite and positive semi-definite. This will happen as long as the expected information per observation

$$E \left\{ \frac{\partial \log L(\theta; Y_i)}{\partial \beta} \frac{\partial \log L(\theta; Y_i)}{\partial \beta'} \right\}$$

is always strictly positive and bounded.

7.5 Generalized Least Squares and instrumental variables**7.5.1 Generalized Least Squares, GLS**

Since Ω is p.d., the inverse exists and is also p.d. Therefore, there exists a nonsingular matrix P so that

$$\Omega^{-1} = P'P.$$

This gives

$$\hat{\beta} = (X'P'PX)^{-1} X'P'Py = \{(PX)'(PX)\}^{-1} (PX)'(Py).$$

This would have been obtained from the regression of Py on PX . To deal with the nonspherical model, pre-multiply $y = X\beta + u$ by P so that

$$y_* = X_*\beta + u_*$$

where $y_* = Py$, $X_* = PX$ and $u_* = Pu$. Since, $\Omega = P^{-1}(P')^{-1}$,

$$\text{var}(u_*) = E(Puu'P') = \sigma^2 P\Omega P' = \sigma^2 I,$$

satisfying the conditions under which OLS is BLUE. The GLS estimator is defined as

$$\beta_{\text{GLS}} = (X_*'X_*)^{-1} X_*'y_* = (X'\Omega^{-1}X)^{-1} (X'\Omega^{-1}y)$$

so that

$$\text{var}(\beta_{\text{GLS}}) = \sigma^2 (X'\Omega^{-1}X)^{-1}$$

which is also the asymptotic variance matrix had we adopted the ML approach. An unbiased estimate of σ^2 can be obtained from OLS applied to the transformed model

$$\begin{aligned} s^2 &= (y_* - X_*\beta_{\text{GLS}})'(y_* - X_*\beta_{\text{GLS}})/(n-k) \\ &= \{P(y - X\beta_{\text{GLS}})\}'\{P(y - X\beta_{\text{GLS}})\}/(n-k) \\ &= (y - X\beta_{\text{GLS}})'\Omega^{-1}(y - X\beta_{\text{GLS}})/(n-k), \end{aligned}$$

which differs from the biased MLE $\hat{\sigma}^2$ by the factor $n/(n-k)$. To test the finite sample restriction $H_0 : R\beta = r$, the test can be based on

$$F = \frac{(r - R\beta_{\text{GLS}})'\{R(X'\Omega^{-1}X)^{-1}R'\}^{-1}(r - R\beta_{\text{GLS}})/q}{s^2} \sim F_{q, n-k}.$$

Note that we could have also written u as $u \sim N(0, V)$ (where $V = \sigma^2\Omega$), so that,

$$\begin{aligned} \beta_{\text{GLS}} &= (X'V^{-1}X)^{-1}X'V^{-1}y \\ \text{var}(\beta_{\text{GLS}}) &= (X'V^{-1}X)^{-1}. \end{aligned}$$

7.5.2 Instrumental Variable (IV) estimators

If the condition stating that the regressors and the disturbance are independent does not hold, then the OLS estimators are biased and inconsistent. As an example, consider the *errors in variables* problem,

$$y = x\beta + u$$

without a constant term, but that the variables are measured with error, that is,

$$\begin{aligned} x &= \bar{x} + v \\ \text{so that } y &= \beta\bar{x} + u \end{aligned}$$

where u , \bar{x} and v are mutually independent. This scenario is classic under models with *measurement error*. The OLS equations give

$$\begin{aligned} \hat{\beta} &= \frac{\sum y_i x_i}{\sum x_i^2} = \frac{\sum x_i(\beta\bar{x}_i + u_i)}{\sum x_i^2} \\ &= \beta \frac{\sum x_i \bar{x}_i}{\sum x_i^2} + \frac{\sum x_i u_i}{\sum x_i^2}. \end{aligned}$$

If the second moments and their plims's exist, then,

$$\begin{aligned} \text{plim} \left(\frac{1}{n} \sum x_i^2 \right) &= \sigma_{\bar{x}}^2 + \sigma_v^2 \\ \text{plim} \left(\frac{1}{n} \sum x_i \bar{x}_i \right) &= \sigma_{\bar{x}}^2 \\ \text{plim} \left(\frac{1}{n} \sum x_i u_i \right) &= 0 \end{aligned}$$

giving

$$\text{plim} \hat{\beta} = \beta \left(\frac{\sigma_{\bar{x}}^2}{\sigma_{\bar{x}}^2 + \sigma_v^2} \right).$$

Hence OLS is biased and inconsistent. Suppose we now re-write the model as

$$y = \beta x + (u - \beta v)$$

so that

$$\hat{\beta} = \beta + \frac{\sum x_i(u_i - \beta v_i)}{\sum x_i^2},$$

then, $\text{plim}(1/n) \sum x_i(u_i - \beta v_i) = \text{plim}(1/n) \sum x_i u_i - \beta \text{plim}(1/n) \sum x_i v_i = -\beta \sigma_v^2$, giving the same results for $\text{plim} \hat{\beta}$ as before. For the general linear model,

$$\hat{\beta} = \beta + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{u}$$

so that,

$$\text{plim } \hat{\beta} = \beta + \text{plim} \left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \cdot \text{plim} \left(\frac{1}{n} \mathbf{X}'\mathbf{u} \right).$$

If we assume that $\text{plim}(\mathbf{X}'\mathbf{X}/n) = \Sigma_{XX}$, a p.d. matrix of full rank, and $\text{plim}(\mathbf{X}'\mathbf{u}/n) = \Sigma_{Xu} \neq 0$, then,

$$\text{plim } \hat{\beta} = \beta + \Sigma_{XX}^{-1} \cdot \Sigma_{Xu}$$

so that correlation of \mathbf{u} with one or more of the disturbance terms makes OLS inconsistent. A consistent estimator can be obtained by using *instrumental variables*.

Suppose the model is still $\mathbf{y} = \mathbf{X}\beta + \mathbf{u}$, with $\text{var}(\mathbf{u}) = \sigma^2 \mathbf{I}$, but $\text{plim}(\mathbf{X}'\mathbf{u}/n) \neq 0$. If we can find a matrix \mathbf{Z} of order $n \times l$, ($l \geq k$), such that

- the variables in \mathbf{Z} are correlated with those in \mathbf{X} and in the limit, $\text{plim}(\mathbf{Z}'\mathbf{X}/n) = \Sigma_{ZX}$, a finite matrix of full rank.
- and $\text{plim}(\mathbf{Z}'\mathbf{u}/n) = 0$.

Then

$$\mathbf{Z}'\mathbf{y} = \mathbf{Z}'\mathbf{X}\beta + \mathbf{Z}'\mathbf{u} \quad \text{with} \quad \text{var}(\mathbf{Z}'\mathbf{u}) = \sigma^2(\mathbf{Z}'\mathbf{Z})$$

which suggests the use of GLS. The estimator is thus

$$\hat{\beta}_{\text{GLS}} = \hat{\beta}_{\text{IV}} = \{\mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}\}^{-1} \mathbf{X}'\mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{y}$$

so that $\hat{\beta}_{\text{GLS}} = \hat{\beta}_{\text{IV}} = (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1} \mathbf{X}'\mathbf{P}_Z\mathbf{y}$, where $\mathbf{P}_Z = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$. The variance-covariance matrix is

$$\text{var}(\hat{\beta}_{\text{IV}}) = \sigma^2(\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1}$$

and σ^2 can be consistently estimated by

$$\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\beta}_{\text{IV}})'(\mathbf{y} - \mathbf{X}\hat{\beta}_{\text{IV}})/n.$$

To check that the IV estimator is consistent,

$$\begin{aligned} \hat{\beta}_{\text{IV}} &= \beta + \left(\frac{1}{n} \mathbf{X}'\mathbf{P}_Z\mathbf{X} \right)^{-1} \left(\frac{1}{n} \mathbf{X}'\mathbf{P}_Z\mathbf{u} \right) \\ \text{and } \frac{1}{n} \mathbf{X}'\mathbf{P}_Z\mathbf{X} &= \left(\frac{1}{n} \mathbf{X}'\mathbf{Z} \right) \left(\frac{1}{n} \mathbf{Z}'\mathbf{Z} \right)^{-1} \left(\frac{1}{n} \mathbf{Z}'\mathbf{X} \right), \end{aligned}$$

so that

$$\text{plim} \left(\frac{1}{n} \mathbf{X}'\mathbf{P}_Z\mathbf{X} \right) = \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZX}$$

which is finite and nonsingular. Also,

$$\text{plim} \left(\frac{1}{n} \mathbf{X}'\mathbf{P}_Z\mathbf{u} \right) = \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{Zu} = 0$$

(since the instruments are assumed to be uncorrelated in the limit with the disturbance), so that the IV estimator is consistent.

Example - special case

If $l = k$, then $\mathbf{X}'\mathbf{Z}$ is a $k \times k$ nonsingular matrix. The estimator now simplifies to

$$\begin{aligned} \hat{\beta}_{\text{IV}} &= (\mathbf{Z}'\mathbf{X})^{-1} \mathbf{Z}'\mathbf{y} \\ \text{with } \text{var}(\hat{\beta}_{\text{IV}}) &= \sigma^2(\mathbf{Z}'\mathbf{X})^{-1}(\mathbf{Z}'\mathbf{Z})(\mathbf{X}'\mathbf{Z})^{-1}. \end{aligned}$$

7.5.2.1 Two-Stage Least Squares (2SLS)

The IV estimator can also be seen as the result of applying least squares twice.

- (1) Regress each of the variables in the \mathbf{X} matrix on \mathbf{Z} to obtain a matrix of fitted values, $\hat{\mathbf{X}}$.

$$\hat{\mathbf{X}} = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X} = \mathbf{P}_Z\mathbf{X}.$$

- (2) Regress \mathbf{y} on $\hat{\mathbf{X}}$ to obtain the estimated β vector

$$\begin{aligned} \hat{\beta}_{\text{2SLS}} &= (\hat{\mathbf{X}}'\hat{\mathbf{X}})^{-1}(\hat{\mathbf{X}}'\mathbf{y}) \\ &= (\mathbf{X}'\mathbf{P}_Z\mathbf{X})^{-1}(\mathbf{X}'\mathbf{P}_Z\mathbf{y}) \\ &= \hat{\beta}_{\text{IV}}. \end{aligned}$$