

# The Incidental Parameters Problem

Bryan S. Graham, UC - Berkeley & NBER

March 26, 2021

In 1948 Jerzy Neyman and Elizabeth Scott published a remarkable paper in *Econometrica*. In fact it was the only paper appearing the January 1948 issues of *Econometrica*. It also represents the only time Neyman, who along with Sir Ronald Fisher, is *the* giant of 20th century statistics, published in an economics journal. Neyman's coauthor, Elizabeth Scott, was a Ph.D. student in Astronomy when the paper was written (and the questions posed in the paper appears to have arisen from questions related to measurement error in astronomical data). Although Scott was a trained astronomer, she was unable to pursue a career in that field because women were prevented from using the telescope at Berkeley at the time. She instead pursued a (distinguished) career in statistics – obtaining an appointment in the Berkeley Math Department in 1951. An award in her honor is presented every two years at the Joint Statistical Meetings.

The following example appears in the paper. Let  $(Y_{i1}, Y_{i2})$  be independent  $\mathcal{N}(\alpha_i, \sigma^2)$  random variables for  $i = 1, \dots, N$ . The location parameter for each  $i$  –  $\alpha_i$  – is different, while the scale parameter,  $\sigma^2$ , is common to all units. We will call  $\{\alpha_i\}_{i=1}^N$  the *incidental parameters* and  $\sigma^2$  the common parameter. It is a good exercise to verify that the MLEs for this problem are

$$\begin{aligned}\hat{\alpha}_i &= \frac{Y_{i1} + Y_{i2}}{2}, \quad i = 1, \dots, N \\ \hat{\sigma}^2 &= \frac{1}{N} \sum_{i=1}^N \frac{(Y_{i1} - \hat{\alpha}_i)^2 + (Y_{i2} - \hat{\alpha}_i)^2}{2} \\ &= \frac{1}{4N} \sum_{i=1}^N (Y_{i1} - Y_{i2})^2.\end{aligned}$$

First observe that  $\hat{\alpha}_i \sim \mathcal{N}\left(\alpha_i, \frac{\sigma^2}{2}\right)$ . The MLEs of the incidental parameters, while in this case unbiased, are not consistent (for  $N \rightarrow \infty$ ). Second, consider the common parameter  $\sigma^2$ . Since,  $\hat{\alpha}_i = \frac{Y_{i1} + Y_{i2}}{2}$  we have  $\hat{\sigma}^2 = \frac{1}{4N} \sum_{i=1}^N (Y_{i1} - Y_{i2})^2$ . Next observe that  $\mathbb{E}[(Y_{i1} - Y_{i2})^2] =$

$2\sigma^2$  and hence  $\hat{\sigma}^2 \xrightarrow{p} \frac{\sigma^2}{2}$ . The estimate of the common parameter is also inconsistent. Like James & Stein (1961) over a decade later, Neyman & Scott (1948) provided an example where MLE performs poorly. This was very much against prevailing intuitions at the time. This example is non-standard because the dimension of the nuisance (incidental) parameter  $\{\alpha_i\}_{i=1}^N$  grows with  $N$ . Classical theorems on the asymptotic properties of MLE do not apply. Settings like this arise naturally in the context of panel data. In that context the incidental parameters model heterogeneity across cross sectional units, while the common parameters are of primary interest (although sometimes we may also be interested in the incidental parameters, or averages of (functions of) them). In this note we'll present a classic example of incidental parameters bias in a non-linear panel data model. The example is both interesting and historically important, but also helpful for developing intuitions about how MLE might behave in more complicated settings with incidental parameters. We will then consider two classic "solutions" to incidental parameters bias: (i) conditioning on a sufficient statistic for them, and (ii) integrated likelihood. The latter approach has natural connections to Bayesian statistics. The survey paper by Arellano (2003) is worth reading. Lancaster (2000) provides a historical survey of the incidental parameters problem. A classic reference in economics is Chamberlain (1980).

## Joint fixed effects maximum likelihood

For each of  $i = 1, \dots, N$  randomly sampled agents we observe a binary choice,  $Y_{it} \in \{0, 1\}$ , and  $K \times 1$  regressor vector,  $X_{it}$ , in each of  $t = 1, \dots, T$  periods. Choice in each period is determined according to the rule

$$Y_{it} = \mathbf{1}(X'_{it}\beta_0 + A_i - U_{it} \geq 0) \quad (1)$$

with  $U_{it}$  an unobserved (standard) logistic random variable; independently and identically distributed across agents and over time. Here  $A_i$  is time-invariant, agent specific heterogeneity that is unobserved by the econometrician.

Let  $X_i = (X'_{i1}, \dots, X'_{iT})'$ ,  $\mathbf{X} = (X_1, \dots, X_N)'$  and  $\mathbf{A} = (A_1, \dots, A_N)'$ . We will begin by considering the properties of the maximum likelihood estimates of  $\beta$  and  $\mathbf{A}$ . We'll call this the joint fixed effects MLE.

The log-likelihood function conditional on  $\mathbf{X}$  is

$$l_N(\beta, \mathbf{A}) = \sum_{i=1}^N l_i(\beta, A_i),$$

where

$$l_i(\beta, A_i) = \sum_{t=1}^T Y_{it} \ln F(X'_{it}\beta + A_i) + (1 - Y_{it}) \ln [1 - F(X'_{it}\beta + A_i)] \quad (2)$$

with

$$F(X'_{it}\beta + A_i) = \frac{\exp(X'_{it}\beta + A_i)}{1 + \exp(X'_{it}\beta + A_i)}.$$

The unit  $i$ 's contribution to the score vectors for the common parameter,  $\beta$ , and the  $N$  incidental parameters,  $\{A_i\}_{i=1}^N$ , are

$$\begin{aligned} S_{\beta i}(\beta, A_i) &= \sum_{t=1}^T (Y_{it} - F(X'_{it}\beta + A_i)) X_{it} \\ S_{A_i}(\beta, A_i) &= \sum_{t=1}^T (Y_{it} - F(X'_{it}\beta + A_i)), \quad i = 1, \dots, N. \end{aligned}$$

Conditional on the common parameter,  $\beta$ , the MLE of  $A_i$  is

$$\hat{A}_i(\beta) = \arg \max_{a \in \mathbb{R}^1} l_i(\beta, a)$$

or the solution to

$$S_{A_i}(\beta, \hat{A}_i(\beta)) = \sum_{t=1}^T (Y_{it} - F(X'_{it}\beta + \hat{A}_i(\beta))) = 0. \quad (3)$$

An important feature of (3) is that, conditional on  $\beta$ , only unit  $i$ 's observations contribute to our estimate of  $A_i$ .

## A special case

In what follows we will consider the case with  $T = 2$ ,  $X_{i1} = (1, 0)'$  and  $X_{i2} = (1, 1)$  for all  $i = 1, \dots, N$  and  $\beta = (\alpha, \gamma)'$ . That is the model consists of an “intercept” and a period 2 time dummy. This is a classic special case considered by Chamberlain (2010) and others. For a given value of  $\beta$ , there will only be four possible solutions to (3); one for each of the four possible sequences for the outcome across the two periods (note that there is no cross-sectional variation in  $X_i$ ).

$(Y_1, Y_2)$	$l_i(\beta, A_i)$	$\hat{A}_i(\beta)$
(0, 0)	$\ln[F(-\alpha - A_i)] + \ln[F(-\alpha - \gamma - A_i)]$	$-\infty$
(0, 1)	$\ln[F(-\alpha - A_i)] + \ln[F(\alpha + \gamma + A_i)]$	$-\alpha - \frac{\gamma}{2}$
(1, 0)	$\ln[F(\alpha + A_i)] + \ln[F(-\alpha - \gamma - A_i)]$	$-\alpha - \frac{\gamma}{2}$
(1, 1)	$\ln[F(\alpha + A_i)] + \ln[F(\alpha + \gamma + A_i)]$	$\infty$

Column 2 of the above table reports the form of the log-likelihood contribution for each of the four possible outcome sequences given in Column 1. You should verify that you can re-construct Column 2 on your own using (2) above. Column 3 reports the MLE of  $A_i$  (conditional on a given value of the common parameter  $\beta$ ).

Consider the (0, 1) case first. Differentiating  $l_i(\beta, A_i)$  with respect to  $A_i$  yields, for the (0, 1) case:

$$\begin{aligned} -\frac{F'(-\alpha - A_i)}{F(-\alpha - A_i)} + \frac{F'(\alpha + \gamma + A_i)}{F(\alpha + \gamma + A_i)} &= -(1 - F(-\alpha - A_i)) + 1 - F(\alpha + \gamma + A_i) \\ &= F(-\alpha - A_i) - F(\alpha + \gamma + A_i). \end{aligned}$$

Therefore  $\hat{A}_i(\beta)$  satisfies

$$\begin{aligned} F(-\alpha - \hat{A}_i(\beta)) &= F(\alpha + \gamma + \hat{A}_i(\beta)) \\ -\alpha - \hat{A}_i(\beta) &= \alpha + \gamma + \hat{A}_i(\beta) \end{aligned}$$

and hence equals

$$\hat{A}_i(\beta) = -\alpha - \frac{\gamma}{2}.$$

For the (1, 0) case we get the same solution. For the (0, 0), the MLE of  $A_i$  is  $-\infty$ ; whereas for the (1, 1) case it is  $\infty$ .

With these four possible values of  $\hat{A}_i(\beta)$  we can form the concentrated log-likelihood

$$l_N^c(\beta) = \sum_{i=1}^N l_i(\beta, \hat{A}_i(\beta)).$$

The first thing to note is that the (0, 0) and (1, 1) cases make zero contribution to the concentrated log-likelihood. In contrast, each (0, 1) unit contributes

$$2 \ln \left[ F\left(\frac{\gamma}{2}\right) \right]$$

and each  $(1, 0)$  unit contributes

$$2 \ln \left[ 1 - F \left( \frac{\gamma}{2} \right) \right].$$

Hence the concentrated log-likelihood equals

$$l_N^c(\beta) = 2 \sum_{i=1}^N \left\{ (1 - Y_{i1}) Y_{i2} \ln \left[ F \left( \frac{\gamma}{2} \right) \right] + Y_{i1} (1 - Y_{i2}) \ln \left[ 1 - F \left( \frac{\gamma}{2} \right) \right] \right\}.$$

Let  $p_0 = \Pr(Y_{i1} = 0, Y_{i2} = 1 | Y_{i1} + Y_{i2} = 1)$ . The MLE of  $p_0$  is

$$\begin{aligned} \hat{p} &= \frac{\sum_{i=1}^N (1 - Y_{i1}) Y_{i2}}{\sum_{i=1}^N \{(1 - Y_{i1}) Y_{i2} + Y_{i1} (1 - Y_{i2})\}} \\ &= \frac{\sum_{i=1}^N \mathbf{1}(Y_{i1} = 0) \mathbf{1}(Y_{i2} = 1)}{\sum_{i=1}^N \mathbf{1}(Y_{i1} + Y_{i2} = 1)} \end{aligned}$$

By the invariance property of MLE we have that

$$\begin{aligned} \hat{\gamma} &= 2F^{-1}(\hat{p}) \\ &\xrightarrow{p} 2F^{-1}(p_0). \end{aligned}$$

Next observe that, using the form of the logit choice model,

$$\begin{aligned} &\Pr(Y_{i1} = 0, Y_{i2} = 1 | Y_{i1} + Y_{i2} = 1, A = a) \\ &= \frac{[1 - F(\alpha_0 + a)] F(\alpha_0 + \gamma_0 + a)}{[1 - F(\alpha_0 + a)] F(\alpha_0 + \gamma_0 + a) + F(\alpha_0 + a) [1 - F(\alpha_0 + \gamma_0 + a)]} \\ &= \frac{\exp(\alpha_0 + \gamma_0 + a)}{\exp(\alpha_0 + \gamma_0 + a) + \exp(\alpha_0 + a)} \\ &= \frac{\exp(\gamma_0)}{1 + \exp(\gamma_0)} \\ &= F(\gamma_0), \end{aligned}$$

The law-of-iterated expectations then gives

$$\begin{aligned} p_0 &= \Pr(Y_{i1} = 0, Y_{i2} = 1 | Y_{i1} + Y_{i2} = 1) \\ &= \mathbb{E}[\Pr(Y_{i1} = 0, Y_{i2} = 1 | Y_{i1} + Y_{i2} = 1, A_i)] \\ &= F(\gamma_0), \end{aligned}$$

and hence that

$$\hat{\gamma} \xrightarrow{p} 2F^{-1}(F(\gamma_0)) = 2\gamma_0.$$

The joint fixed MLE of the common parameter  $\gamma$  is inconsistent. This is another example of incidental parameters bias.

## Conditional fixed effects maximum likelihood

One intuition for the poor properties of the joint fixed effects MLEs is that the badly estimated incidental parameters “contaminate” our estimates of the common parameter. Conditional fixed effects approaches search for implications of the model which are invariant to the incidental parameters, but which still contain identifying information about  $\beta$ . For a certain class of models with exponential family structure, such estimates are available as conditional maximum likelihood estimators. The conditional fixed effect logit estimator is one such estimator.

Re-arranging the log-likelihood yields

$$\begin{aligned}
 L(\beta, \mathbf{A}) &= \Pr(\mathbf{Y} = \mathbf{y} | \mathbf{X}, \mathbf{A}) \\
 &= \prod_{i=1}^N \prod_{t=1}^T \left[ \frac{\exp(X'_{it}\beta + A_i)}{1 + \exp(X'_{it}\beta + A_i)} \right]^{Y_{it}} \left[ \frac{1}{1 + \exp(X'_{it}\beta + A_i)} \right]^{1-Y_{it}} \\
 &= \prod_{i=1}^N \prod_{t=1}^T \left[ \frac{\exp(X'_{it}\beta + A_i)^{Y_{it}}}{1 + \exp(X'_{it}\beta + A_i)} \right] \\
 &= c(\mathbf{X}; \beta, \mathbf{A}) \exp \left( \left[ \sum_{i=1}^N \sum_{t=1}^T Y_{it} X'_{it} \right] \beta + \left[ \sum_{t=1}^T Y_{1t} \right] A_1 + \cdots + \left[ \sum_{t=1}^T Y_{Nt} \right] A_N \right) \\
 &= c(\mathbf{X}; \beta, \mathbf{A}) \prod_{i=1}^N \exp \left( \sum_{t=1}^T Y_{it} X'_{it} \beta + \left[ \sum_{t=1}^T Y_{it} \right] A_i \right).
 \end{aligned}$$

Hence  $\sum_{t=1}^T Y_{it}$  is a sufficient statistic for  $A_i$  and we have that

$$\begin{aligned}
 \Pr \left( Y_i = y | X_i, A_i, \sum_{t=1}^T Y_{it} \right) &= \frac{\exp \left( \sum_{t=1}^T Y_{it} X'_{it} \beta \right)}{\sum_{v \in \mathbb{V}_i} \exp \left( \sum_{t=1}^T v_t X'_{it} \beta \right)} \\
 &= \Pr \left( Y_i = y | X_i, \sum_{t=1}^T Y_{it} \right)
 \end{aligned}$$

where

$$\mathbb{V}_i = \left\{ v : v \in \{0, 1\}^T, \sum_{t=1}^T v_t = \sum_{t=1}^T Y_{it} \right\}.$$

Conditional on  $\sum_{t=1}^T Y_{it}$ , the likelihood of the event  $Y_{i1} = y_1, \dots, Y_{iT} = y_T$  does not vary with  $A_i$ . Hence we can choose  $\hat{\beta}$  to maximize

$$l_N^{\text{cond}}(\beta) = \sum_{i=1}^N \left\{ \left[ \sum_{t=1}^T Y_{it} X'_{it} \right] \beta - \ln \left[ \sum_{v \in \mathbb{V}_i} \exp \left( \sum_{t=1}^T v_t X'_{it} \beta \right) \right] \right\}.$$

This yields the conditional fixed effects maximum likelihood estimate of  $\beta$ . The main challenge in computing  $\hat{\beta}$  involves efficiently evaluating  $\sum_{v \in \mathbb{V}_i} \exp \left( \sum_{t=1}^T v_t X'_{it} \beta \right)$ . Historically this was actually a serious problem. Fortunately at some point it was realized that this denominator could be computed using a recursive algorithm.

Let  $S_i = \sum_{t=1}^T Y_{it}$  be the number of times unit  $i$ 's outcome equals 1. The set  $\mathbb{V}_i$  includes  $\binom{T}{S_i}$  elements. This set consists of all possible sequences  $Y_{i1}, \dots, Y_{iT}$  such that  $S_i = \sum_{t=1}^T Y_{it}$  (think in terms of Bernoulli trials). Next define

$$\begin{aligned} w_i(T, S_i) &= \sum_{v \in \mathbb{V}_i} \exp \left( \sum_{t=1}^T v_t X'_{it} \beta \right) \\ &= \sum_{v \in \mathbb{V}_i} \exp \left( \sum_{t=1}^{T-1} v_t X'_{it} \beta \right) \exp(v_T X'_{iT} \beta) \end{aligned}$$

which equals

$$w_i(T-1, S_i)$$

if  $v_T = 0$  and equals

$$w_i(T-1, S_i-1) \exp(X'_{iT} \beta)$$

if  $v_T = 1$ . We therefore have the recursive relationship

$$w_i(T, S_i) = w_i(T-1, S_i) + w_i(T-1, S_i-1) \exp(X'_{iT} \beta).$$

The first term above captures the contribution of all sequences in  $\mathbb{V}_i$  which end with  $v_T = 0$ , whereas the second term is the contribution of all those sequences with  $v_T = 1$ . We set  $w_i(T, s) = 0$  if  $T < s$  and  $w_i(T, 0) = 1$ . The denominator of the conditional likelihood can be built up recursively using relationships like the one sketched above.

When  $T = 2$  the conditional logit estimator takes a particular simple form. In that case we have

$$\begin{aligned}
\Pr(Y_{i1} = 1 | X_i, A_i, Y_{i1} + Y_{i2} = 1) &= \frac{\Pr(Y_{i1} = 1, Y_{i2} = 0 | X_i, A_i)}{\Pr(Y_{i1} = 1, Y_{i2} = 0 | X_i, A_i) + \Pr(Y_{i1} = 0, Y_{i2} = 1 | X_i, A_i)} \\
&= \frac{\frac{\exp(X'_{i1}\beta + A_i)}{1 + \exp(X'_{i1}\beta + A_i)} \frac{1}{1 + \exp(X'_{i2}\beta + A_i)}}{\frac{\exp(X'_{i1}\beta + A_i)}{1 + \exp(X'_{i1}\beta + A_i)} \frac{1}{1 + \exp(X'_{i2}\beta + A_i)} + \frac{1}{1 + \exp(X'_{i1}\beta + A_i)} \frac{\exp(X'_{i2}\beta + A_i)}{1 + \exp(X'_{i2}\beta + A_i)}} \\
&= \frac{\exp(X'_{i1}\beta + A_i)}{\exp(X'_{i1}\beta + A_i) + \exp(X'_{i2}\beta + A_i)} \\
&= \frac{1}{1 + \exp((X_{i2} - X_{i1})' \beta)}.
\end{aligned}$$

Similar calculations give  $\Pr(Y_{i1} = 0 | X_i, A_i, Y_{i1} + Y_{i2} = 1) = \frac{\exp((X_{i2} - X_{i1})' \beta)}{1 + \exp((X_{i2} - X_{i1})' \beta)}$ . Hence, in the  $T = 2$  case, the conditional fixed effects MLE maximizes

$$l_N^{\text{cond}}(\beta) = 2 \sum_{i=1}^N \left\{ (1 - Y_{i1}) Y_{i2} \ln \left[ \frac{\exp((X_{i2} - X_{i1})' \beta)}{1 + \exp((X_{i2} - X_{i1})' \beta)} \right] + Y_{i1} (1 - Y_{i2}) \ln \left[ \frac{1}{1 + \exp((X_{i2} - X_{i1})' \beta)} \right] \right\}$$

## Integrated likelihood or random effects approaches

An alternative, and indeed quite natural, approach to large numbers of incidental parameters is to assign a prior distribution to them. One can then integrate these parameters out of the likelihood. The resultant criterion function is an integrated, random effects or empirical Bayes likelihood. I will have more to say about this approach in the next lecture. A quick sketch starts with the same choice model

$$Y_{it} = \mathbf{1}(X'_{it}\beta + A_i - V_{it} \geq 0).$$

However now it is convenient to make a normal assumption on the time-varying shocks

$$V_{it} | X_{i1}, \dots, X_{iT}, A_i \sim \mathcal{N}(0, 1).$$

Our earlier approaches left the joint distribution of  $A_i$  and  $X_{i1}, \dots, X_{iT}$  unrestricted. In (correlated) random effects approaches we will restrict this relationship. Here we will assume that

$$A_i | X_{i1}, \dots, X_{iT} \sim \mathcal{N}(X'_i \pi, \sigma_A^2).$$

Substitution yields

$$Y_{it} = \mathbf{1}(X'_{it}\beta + X'_i \pi + A_i^* \geq V_{it})$$



such that

$$\Pr(Y_{it} = 1 | X_i, A_i^*) = \Phi(X'_{it}\beta + X'_i\pi + A_i^*).$$

The correlated random effects estimator choose  $\beta$  and  $\eta = (\pi', \sigma_A)'$  to maximize

$$l^I(Y_i | X_i; \beta, \eta) = \ln \left[ \int \prod_{t=1}^T \Phi(X'_{it}\beta + X'_i\pi + a)^{Y_{it}} [1 - \Phi(X'_{it}\beta + X'_i\pi + a)]^{1-Y_{it}} \frac{1}{\sigma_A} \phi\left(\frac{a}{\sigma_A}\right) da \right].$$

## References

- Arellano, M. (2003). Discrete choices with panel data. *Investigaciones Económicas*, 27(3), 423 – 458.
- Chamberlain, G. (1980). Analysis of covariance with qualitative data. *Review of Economic Studies*, 47(1), 225 – 238.
- Chamberlain, G. (2010). Binary response models for panel data: Identification and information. *Econometrica*, 78(1), 159 – 168.
- James, W. & Stein, C. M. (1961). Estimation with quadratic loss. *Berkeley Symposium on Mathematical Statistics and Probability*, 4(1), 361 – 379.
- Lancaster, T. (2000). The incidental parameter problem since 1948. *Journal of Econometrics*, 95(2), 391 – 413.
- Neyman, J. & Scott, E. L. (1948). Consistent estimates based on partially consistent observations. *Econometrica*, 16(1), 1 – 32.