

# Model Based Instrumental Variables

Bryan S. Graham, UC - Berkeley & NBER

November 26, 2019

Angrist (1990) studies the relationship between military service in Vietnam and earnings for American men born in the early 1950s. For concreteness let the population of interest be American men born in 1950. Let  $Y$  denote earnings in 1981 for a random draw from this population; let  $D$  equal one if this individual served in Vietnam and zero otherwise; finally, let  $X$  equal one if the individual was “draft eligible” (i.e., had a low draft lottery number) and zero otherwise.

Draft eligibility ( $X = 1$ ) encourages military service. In what follows we will refer to  $X$  as the *encouragement* or *instrumental variable*. Let  $D(1) \in \{0, 1\}$  denote an individual’s veteran status when, possibly contrary to fact, they are draft eligible. Let  $D(0) \in \{0, 1\}$  denote their veteran status, again possibly contrary to fact, when they are draft ineligible. We can partition individuals into four strata based upon their military service behavior with and without encouragement.

Table 1 depicts the four possible *compliance strata*. First, there are *never-takers* ( $D(0) = 0, D(1) = 0$ ), who never serve in the military irrespective of draft eligibility. Second, there are *always-takers* ( $D(0) = 1, D(1) = 1$ ), who serve in the military regardless of draft eligibility. Third, and important for our analysis, are *compliers* ( $D(0) = 0, D(1) = 1$ ), who serve when encouraged (i.e., when draft eligible), but who do not serve when not encouraged. Finally, there are *defiers* ( $D(0) = 1, D(1) = 0$ ), who avoid service when encouraged, but – inexplicably – serve when not. Will be rule out the existence of this last compliance stratum by making the **monotonicity assumption** that  $D(1) \geq D(0)$ : no unit is less likely to serve when encouraged to do so.

Table 1: Compliance Strata

	$D(1) = 0$	$D(1) = 1$
$D(0) = 0$	Never-Taker	Complier
$D(0) = 1$	Defier	Always-Taker

Let  $Y(d, x)$  denote an individual's earnings in 1981 when, possibly contrary to fact, their military service and draft eligibility are, respectively,  $D = d$  and  $X = x$ . We will additionally make the **exclusion restriction** that  $Y(d, 1) = Y(d, 0)$  for  $d = 0, 1$ . Conditional on serving in the military, whether one was draft eligible or not has no effect on earnings. Likewise conditional on not serving in the military, draft eligibility has no effect on earning. An implication of this assumption is that any effect of draft eligibility on earnings must operate by inducing individuals to serve or not to serve.

Our final assumption is that the encouragement/instrument is **randomly assigned**.

$$(D(0), D(1), Y(0), Y(1)) \perp X.$$

The causal effect of military service on earnings for individual  $i$  equals

$$Y_i(1) - Y_i(0).$$

While this effect is well defined, it is never observable since we either observe earnings given military service or not, but never both simultaneously. Indeed observed earnings equal

$$Y = (1 - D(X))Y(0) + D(X)Y(1). \tag{1}$$

While individual-level treatment effects are never identified, our hope is that average effects are.

## Identification

Let  $A$  be a  $3 \times 1$  vector of compliance strata indicator variables. If an individual is a never-taker the first element of this vector is a one, with the other elements being zero. If he is a complier, the second element is a one, with the other elements being zero. Finally, if he is an always-taker the last element is a one with the balance being zero. Recall that there are no defiers when we maintain the monotonicity assumption.

Let  $a_j$  be a vector of zeros with the exception of the  $j^{th}$  element, which is equal to one. Further let  $\pi_j = \Pr(A = a_j)$  equal the population frequency of compliance strata  $j$ . By the

law of total probability

$$\begin{aligned}
\Pr(Y \leq y | X = 0) &= \Pr(A = a_1 | X = 0) \Pr(Y \leq y | X = 0, A = a_1) \\
&\quad + \Pr(A = a_2 | X = 0) \Pr(Y \leq y | X = 0, A = a_2) \\
&\quad + \Pr(A = a_3 | X = 0) \Pr(Y \leq y | X = 0, A = a_3) \\
&= \pi_1 \Pr(Y(0) \leq y | A = a_1) + \pi_2 \Pr(Y(0) \leq y | A = a_2) \\
&\quad + \pi_3 \Pr(Y(1) \leq y | A = a_1)
\end{aligned} \tag{2}$$

An individual's compliance strata  $A$  is fully determined by their configuration of  $(D(0), D(1))$ . By random assignment we therefore have that  $A$  and  $X$  are independent (i.e.,  $\Pr(A = a_j | X) = \Pr(A = a_j)$ ). Note also that  $D$  is a degenerate random variable given a unit's compliance strata and encouragement. We also have that  $Y(d)$  is independent of  $X$  conditional on  $A$ . The second equality above follows from these facts, equation (1), and the military service behavior of units in each stratum with and without encouragement.

Analogous arguments give

$$\begin{aligned}
\Pr(Y \leq y | X = 1) &= \pi_1 \Pr(Y(0) \leq y | A = a_1) + \pi_2 \Pr(Y(1) \leq y | A = a_2) \\
&\quad + \pi_3 \Pr(Y(1) \leq y | A = a_1).
\end{aligned} \tag{3}$$

Finally we have (try to show this on your own):

$$\Pr(D = 1 | X = 0) = 0 \cdot \pi_1 + 0 \cdot \pi_2 + 1 \cdot \pi_3 \tag{4}$$

$$\Pr(D = 1 | X = 1) = 0 \cdot \pi_1 + 1 \cdot \pi_2 + 1 \cdot \pi_3. \tag{5}$$

Together (2),(3),(4) and (5) yield

$$\Pr(Y(1) \leq y | A = a_2) - \Pr(Y(0) \leq y | A = a_2) = \frac{\Pr(Y \leq y | X = 1) - \Pr(Y \leq y | X = 0)}{\Pr(D = 1 | X = 1) - \Pr(D = 1 | X = 0)}.$$

The difference in the marginal distribution of the potential outcome under treatment,  $Y(1)$ , versus that under control,  $Y(0)$ , is identified for the subpopulation of compliers (Imbens & Rubin, 1997). A consequence is that the **local average treatment effect** (LATE)

$$\beta^{\text{LATE}} = \mathbb{E}[Y(1) - Y(0) | A = a_2] = \frac{\mathbb{E}[Y | X = 1] - \mathbb{E}[Y | X = 0]}{\Pr(D = 1 | X = 1) - \Pr(D = 1 | X = 0)}$$

is also identified. Since the expression to the right of the equality equals  $\mathbb{C}(Y, X) / \mathbb{C}(D, X)$  the LATE also coincides with the probability limit of the coefficient on  $D$  in the instrumental

variables fit of  $Y$  onto a constant and  $D$  using  $X$  as an instrument for  $D$  (Angrist et al., 1996).

## Likelihood

Our goal is to write down a likelihood for  $(D, Y)$  given  $X$ . To do this we will first write down a likelihood for  $(A, D, Y)$  given  $X$ . We will then treat a unit's compliance strata as missing data. We will call the latter likelihood the **complete data likelihood**, the former the observed or **integrated likelihood**.

We begin by observing that

$$\begin{aligned} f(A, D, Y | X) &= f(D, Y | X, A) f(A | X) \\ &= f(D, Y | X, A) f(A). \end{aligned}$$

The second equality follows by random assignment of  $X$ .

We will assume that the distributions of potential (log) earnings given compliance strata are Gaussian, with strata-specific location and scale parameters (other parametric assumptions could be made here depending on the outcome of interest). We have that the distribution of potential earnings given non-service among never-takers is

$$Y(0) | A = a_1 \sim \mathcal{N}(\mu_{N0}, \sigma_{N0}^2).$$

Since never-takers never serve the distribution of  $Y(1) | A = a_1$  is undefined. For compliers we do need to define both potential earning distributions:

$$\begin{aligned} Y(0) | A = a_2 &\sim \mathcal{N}(\mu_{C0}, \sigma_{C0}^2). \\ Y(1) | A = a_2 &\sim \mathcal{N}(\mu_{C1}, \sigma_{C1}^2). \end{aligned}$$

Finally, for always-takers we need only the  $Y(1)$  distribution:

$$Y(1) | A = a_3 \sim \mathcal{N}(\mu_{A1}, \sigma_{A1}^2).$$

Let  $\phi(y; \mu, \sigma^2)$  be the density of a  $\mathcal{N}(\mu, \sigma^2)$  random variable at  $Y = y$ . With these assumptions we compute conditional densities for service and earnings across the three compliance strata of

$$\begin{aligned}
f(D, Y | X, A = a_1) &= \phi(Y; \mu_{N0}, \sigma_{N0}^2)^{(1-D)(1-X)} \phi(Y; \mu_{N0}, \sigma_{N0}^2)^{(1-D)X} 0^{D(1-X)} 0^{DX} \\
f(D, Y | X, A = a_2) &= \phi(Y; \mu_{C0}, \sigma_{C0}^2)^{(1-D)(1-X)} 0^{(1-D)X} 0^{D(1-X)} \phi(Y; \mu_{C1}, \sigma_{C1}^2)^{DX} \\
f(D, Y | X, A = a_3) &= 0^{(1-D)(1-X)} 0^{(1-D)X} \phi(Y; \mu_{A1}, \sigma_{A1}^2)^{D(1-X)} \phi(Y; \mu_{A1}, \sigma_{A1}^2)^{DX}.
\end{aligned}$$

You should check that the above densities equal zero for unsupported events.

With the above notation we get a complete data log-likelihood contribution for unit  $i$  of

$$\ln L_i^C(\theta, A_i) = \sum_{j=1}^3 A_{ji} \{ \ln f(D_i, Y_i | X_i, A_i = a_j) + \ln \pi_j \}. \quad (6)$$

## Complete data analysis

It is instructive to first consider an analysis which presumes compliance strata are observed. Assume we have access to the random sample  $\{A_i, D_i, Y_i, X_i\}_{i=1}^N$ . Since the compliance strata population shares sum to one we maximize the Lagrangian

$$\mathcal{L} = \sum_{i=1}^N \sum_{j=1}^3 A_{ji} \{ \ln f(D_i, Y_i | X_i, A_i = a_j) + \ln \pi_j \} + \lambda(1 - \pi_1 - \pi_2 - \pi_3)$$

yielding MLEs of

$$\hat{\pi}_1 = \frac{1}{N} \sum A_{1i}, \quad \hat{\pi}_2 = \frac{1}{N} \sum A_{2i}, \quad \hat{\pi}_3 = \frac{1}{N} \sum A_{3i}$$

and

$$\begin{aligned}
\hat{\mu}_{N0} &= \frac{\sum_{i=1}^N A_{1i} \{(1 - D_i)(1 - X_i) + (1 - D_i)X_i\} Y_i}{\sum_{i=1}^N A_{1i} \{(1 - D_i)(1 - X_i) + (1 - D_i)X_i\}}, \\
\hat{\sigma}_{N0}^2 &= \frac{\sum_{i=1}^N A_{1i} \{(1 - D_i)(1 - X_i) + (1 - D_i)X_i\} (Y_i - \hat{\mu}_{N0})^2}{\sum_{i=1}^N A_{1i} \{(1 - D_i)(1 - X_i) + (1 - D_i)X_i\}}, \\
\hat{\mu}_{C0} &= \frac{\sum_{i=1}^N A_{2i} (1 - D_i)(1 - X_i) Y_i}{\sum_{i=1}^N A_{2i} (1 - D_i)(1 - X_i)}, \quad \hat{\sigma}_{C0}^2 = \frac{\sum_{i=1}^N A_{2i} (1 - D_i)(1 - X_i) (Y_i - \hat{\mu}_{C0})^2}{\sum_{i=1}^N A_{2i} (1 - D_i)(1 - X_i)}, \\
\hat{\mu}_{C1} &= \frac{\sum_{i=1}^N A_{2i} D_i X_i Y_i}{\sum_{i=1}^N A_{2i} D_i X_i}, \quad \hat{\sigma}_{C1}^2 = \frac{\sum_{i=1}^N A_{2i} D_i X_i (Y_i - \hat{\mu}_{C1})^2}{\sum_{i=1}^N A_{2i} D_i X_i}, \\
\hat{\mu}_{A1} &= \frac{\sum_{i=1}^N A_{3i} \{D_i(1 - X_i) + D_i X_i\} Y_i}{\sum_{i=1}^N A_{3i} \{D_i(1 - X_i) + D_i X_i\}}, \\
\hat{\sigma}_{A1}^2 &= \frac{\sum_{i=1}^N A_{3i} \{D_i(1 - X_i) + D_i X_i\} (Y_i - \hat{\mu}_{A1})^2}{\sum_{i=1}^N A_{3i} \{D_i(1 - X_i) + D_i X_i\}}.
\end{aligned} \tag{7}$$

Consider the mean outcome for compliers under military service,  $\mu_{C1}$ . If knowledge of the compliance strata were available we would estimate this by the average outcome among compliers who served in the military. Why is this okay? Within a compliance strata the only reason why one unit serves and another doesn't is because those units received different encouragements (in this case a different draft lottery number). Since encouragement is randomly assigned, then so is actual treatment *within a compliance strata*. Encouraged compliers (who serve) and non-encouraged compliers (who do not serve) are fully comparable. Therefore

$$\hat{\beta}^{\text{LATE}} = \hat{\mu}_{C1} - \hat{\mu}_{C0}$$

provides a consistent estimate of the LATE, or the average effect of military service on earnings among compliers. Note we cannot estimate an average effect for never-takers or always-takers, since in these two strata only one of the two potential outcome distributions is identified.

## Incomplete data analysis

In practice compliance strata are unobserved. Let  $\tilde{A}_{1i} \in [0, 1]$  be the subjective probability the econometrician attaches to the statement “unit  $i$  is a never-taker” being true. Let  $\delta = (\mu_{N0}, \sigma_{N0}^2, \mu_{C0}, \sigma_{C0}^2, \mu_{C1}, \sigma_{C1}^2, \mu_{N0}, \sigma_{N0}^2)'$  and  $\pi = (\pi_1, \pi_2, \pi_3)'$ . For the time being assume that the parameter  $\theta = (\delta', \pi')'$  is known. Our likelihood, indexed by  $\theta$ , and the data can be used – via Bayes' rule – to compute the probability that unit  $i$  belongs to a certain

compliance strata.

The probability that “unit  $i$  is a never-taker” given that he did not serve *and* was encouraged to do so is 1:

$$\Pr(A_i = a_1 | D_i = 0, Y_i, X_i = 1; \theta) = 1. \quad (8)$$

This follows because always-takers always serve, and compliers served if encouraged, hence any encouraged unit *not* serving is definitely a never-taker. The situation is a bit more complicated for non-serving units who were not encouraged. An unencouraged unit who did not serve could be either be a complier or a never-taker. To compute the probability of the latter we use Bayes’ rule:

$$\Pr(A_i = a_1 | D_i = 0, Y_i, X_i = 0; \theta) = \frac{f(D_i = 0, Y_i | X_i = 0, A_i = a_1; \theta) \Pr(A_i = a_1 | X_i = 0; \theta)}{f(D_i = 0, Y_i | X_i = 0; \theta)}.$$

We have that  $\Pr(A_i = a_1 | X_i = 0; \theta) = \pi_1$  and  $f(D_i = 0, Y_i | X_i = 0, A_i = a_1; \theta) = \phi(Y_i; \mu_{N0}, \sigma_{N0}^2)$ . This gives the numerator in the expression above.

Since non-serving units who were not encouraged are a mixture of never-takers and compliers, the marginal density  $f(D_i = 0, Y_i | X_i = 0; \theta)$  equals

$$f(D_i = 0, Y_i | X_i = 0; \theta) = \pi_1 \phi(Y_i; \mu_{N0}, \sigma_{N0}^2) + \pi_2 \phi(Y_i; \mu_{C0}, \sigma_{C0}^2).$$

To verify this claim use the identity

$$f(D_i = 0, Y_i | X_i = 0; \theta) = \sum_{j=1}^3 f(D_i = 0, Y_i | X_i = 0, A_i = a_j; \theta) \Pr(A_i = a_j | X_i = 0; \theta)$$

and also note that since always-takers always serve  $f(D_i = 0, Y_i | X_i = 0, A_i = a_3; \theta) = 0$ . Putting these results together yields

$$\Pr(A_i = a_1 | D_i = 0, Y_i, X_i = 0; \theta) = \frac{\pi_1 \phi(Y_i; \mu_{N0}, \sigma_{N0}^2)}{\pi_1 \phi(Y_i; \mu_{N0}, \sigma_{N0}^2) + \pi_2 \phi(Y_i; \mu_{C0}, \sigma_{C0}^2)}. \quad (9)$$

From (8) and (9) we get a posterior probability of unit  $i$  being a never taker of

$$\tilde{A}_{1i} = (1 - X_i)(1 - D_i) \frac{\pi_1 \phi(Y_i; \mu_{N0}, \sigma_{N0}^2)}{\pi_1 \phi(Y_i; \mu_{N0}, \sigma_{N0}^2) + \pi_2 \phi(Y_i; \mu_{C0}, \sigma_{C0}^2)} + X_i(1 - D_i). \quad (10)$$

Note that (10) evaluates to zero whenever  $D_i = 1$  (since never-takers never serve in the military). Similarly if  $X_i = 1$  and  $D_i = 0$  it evaluates to one, since an encouraged unit who does not serve is a never-taker with probability one. The interesting case is when both  $X_i$

and  $D_i$  equal zero; units in this group constitute a mixture of never takers and compliers. Similar arguments yield

$$\begin{aligned} \tilde{A}_{2i} = (1 - X_i) (1 - D_i) & \frac{\pi_2 \phi(Y_i; \mu_{C0}, \sigma_{C0}^2)}{\pi_1 \phi(Y_i; \mu_{N0}, \sigma_{N0}^2) + \pi_2 \phi(Y_i; \mu_{C0}, \sigma_{C0}^2)} \\ & + X_i D_i \frac{\pi_2 \phi(Y_i; \mu_{C1}, \sigma_{C1}^2)}{\pi_2 \phi(Y_i; \mu_{C1}, \sigma_{C1}^2) + \pi_3 \phi(Y_i; \mu_{A1}, \sigma_{A1}^2)}. \end{aligned} \quad (11)$$

and

$$\tilde{A}_{3i} = (1 - X_i) D_i + X_i D_i \frac{\pi_3 \phi(Y_i; \mu_{A1}, \sigma_{A1}^2)}{\pi_2 \phi(Y_i; \mu_{C1}, \sigma_{C1}^2) + \pi_3 \phi(Y_i; \mu_{A1}, \sigma_{A1}^2)}. \quad (12)$$

It is a good exercise to verify that expressions (11) and (12) are correct.

Since  $A_i$  is unobserved we can not evaluate unit  $i$ 's contribution to the complete data log-likelihood (6). However, using (10), (11) and (12), we can evaluate this unit's expected log-likelihood contribution:

$$\mathbb{E} \left[ \ln L_i^C(\theta, A_i) \mid X_i, D_i, Y_i; \hat{\theta}^{(s)} \right] = \sum_{j=1}^3 \tilde{A}_{ji} \{ \ln f(D_i, Y_i \mid X_i, A_i = a_j) + \ln \pi_j \} \quad (13)$$

Note that the value of  $\theta$  used to compute the above expectation, which involves an average over the posterior distribution of  $A_i$ , may differ from the value of  $\theta$  at which the expected log-likelihood is evaluated. We will use  $\hat{\theta}^{(s)}$  to denote the former and  $\theta$  to denote the latter. Summing over all  $i = 1, \dots, N$  units we get a criterion function which is proportion to something called the  $Q$ -function (defined carefully below):

$$Q(\theta, \hat{\theta}^{(s)}) \propto \sum_{i=1}^N \sum_{j=1}^3 \tilde{A}_{ji} \{ \ln f(D_i, Y_i \mid X_i, A_i = a_j) + \ln \pi_j \} \quad (14)$$

Maximizing (14) with respect to  $\theta$  yields

$$\hat{\pi}_1^{(s+1)} = \frac{1}{N} \sum \tilde{A}_{1i}, \quad \hat{\pi}_2^{(s+1)} = \frac{1}{N} \sum \tilde{A}_{2i}, \quad \hat{\pi}_3^{(s+1)} = \frac{1}{N} \sum \tilde{A}_{3i}$$



and

$$\begin{aligned}
\hat{\mu}_{N0}^{(s+1)} &= \frac{\sum_{i=1}^N \tilde{A}_{1i} \{(1-D_i)(1-X_i) + (1-D_i)X_i\} Y_i}{\sum_{i=1}^N \tilde{A}_{1i} \{(1-D_i)(1-X_i) + (1-D_i)X_i\}}, \\
(\hat{\sigma}_{N0}^2)^{(s+1)} &= \frac{\sum_{i=1}^N \tilde{A}_{1i} \{(1-D_i)(1-X_i) + (1-D_i)X_i\} \left(Y_i - \hat{\mu}_{N0}^{(s+1)}\right)^2}{\sum_{i=1}^N \tilde{A}_{1i} \{(1-D_i)(1-X_i) + (1-D_i)X_i\}}, \\
\hat{\mu}_{C0}^{(s+1)} &= \frac{\sum_{i=1}^N \tilde{A}_{2i} (1-D_i)(1-X_i) Y_i}{\sum_{i=1}^N \tilde{A}_{2i} (1-D_i)(1-X_i)}, \quad (\hat{\sigma}_{C0}^2)^{(s+1)} = \frac{\sum_{i=1}^N \tilde{A}_{2i} (1-D_i)(1-X_i) \left(Y_i - \hat{\mu}_{C0}^{(s+1)}\right)^2}{\sum_{i=1}^N \tilde{A}_{2i} (1-D_i)(1-X_i)}, \\
\hat{\mu}_{C1}^{(s+1)} &= \frac{\sum_{i=1}^N \tilde{A}_{2i} D_i X_i Y_i}{\sum_{i=1}^N \tilde{A}_{2i} D_i X_i}, \quad (\hat{\sigma}_{C1}^2)^{(s+1)} = \frac{\sum_{i=1}^N \tilde{A}_{2i} D_i X_i \left(Y_i - \hat{\mu}_{C1}^{(s+1)}\right)^2}{\sum_{i=1}^N \tilde{A}_{2i} D_i X_i}, \\
\hat{\mu}_{A1}^{(s+1)} &= \frac{\sum_{i=1}^N \tilde{A}_{3i} \{D_i(1-X_i) + D_i X_i\} Y_i}{\sum_{i=1}^N \tilde{A}_{3i} \{D_i(1-X_i) + D_i X_i\}}, \\
(\hat{\sigma}_{A1}^2)^{(s+1)} &= \frac{\sum_{i=1}^N \tilde{A}_{3i} \{D_i(1-X_i) + D_i X_i\} \left(Y_i - \hat{\mu}_{A1}^{(s+1)}\right)^2}{\sum_{i=1}^N \tilde{A}_{3i} \{D_i(1-X_i) + D_i X_i\}}.
\end{aligned}$$

These expressions are identical to those derived in the complete data case after replacing stratum indicator  $A_{1i}$  with the posterior probability  $\tilde{A}_{1i}$  and so on.

The above discussion suggests the following estimation algorithm.

1. Let  $\hat{\theta}^{(s)}$  for  $s = 0$  be an initial value for the parameter  $\theta$ .
2. **E-Step:** Compute the expected log-likelihood given the data and current parameter value  $\hat{\theta}^{(s)}$  according to equation (14).
3. **M-Step:** Maximize  $Q(\theta, \hat{\theta}^{(s)})$  with respect to  $\theta$ . Denote the solution to this maximization problem by  $\hat{\theta}^{(s+1)}$ .
4. Repeat steps 2 and 3 until  $Q(\hat{\theta}^{(s+1)}, \hat{\theta}^{(s)}) \approx Q(\hat{\theta}^{(s)}, \hat{\theta}^{(s-1)})$  is small and/or until  $\hat{\theta}^{(s+1)} \approx \hat{\theta}^{(s)}$ .

Both the E- and M-steps have closed-form solutions in the current setting.

## Theory of the EM Algorithm

Unit  $i$ 's contribution to the **observed log-likelihood** is computed by “averaging over” the compliance strata distribution:

$$l_i^I(\theta) = \ln \left( \sum_{l=1}^K L_i^C(\theta, a_l) \right).$$

The integrate log-likelihood for the entire sample is then

$$l_N^I(\theta) = \sum_{i=1}^N l_i^I(\theta). \quad (15)$$

Let  $q(a)$  be some assignment of probability mass to the  $K$  possible types such that  $q(a_k) > 0$  for all  $k = 1, \dots, K$  and  $\sum_{k=1}^K q(a_k) = 1$ . Here “types” correspond to the three compliance strata so that  $K = 3$ .

We can show that the  $i^{th}$  unit's contribution to the observed log likelihood is bounded below by

$$\begin{aligned} \ln \left( \sum_{l=1}^K L_i^C(\theta, a_l) \right) &= \ln \left( \sum_{l=1}^K q(a_l) \frac{L_i^C(\theta, a_l)}{q(a_l)} \right) \\ &\geq \sum_{l=1}^K q(a_l) \ln \left( \frac{L_i^C(\theta, a_l)}{q(a_l)} \right) \\ &= Q_i^*(\theta, q) \end{aligned} \quad (16)$$

where the middle line follows from Jensen's inequality:  $g(\mathbb{E}[Y]) \geq \mathbb{E}[g(Y)]$  for  $g(\cdot)$  concave. Here  $\ln(\cdot)$  is concave and expectations are with respect to  $q(a)$ . The last line defines  $Q_i^*(\theta, q)$ . Equation (16) gives

$$l_N^I(\theta) \geq \sum_{i=1}^N Q_i^*(\theta, q_i)$$

for any set of valid distribution functions  $\{q_i\}_{i=1}^N$  that assign positive probability to each  $\{a_k\}_{k=1}^K$ .

Bayes' Theorem, and the form of the complete data likelihood ((6)), yields the conditional type probabilities

$$\Pr(A = a_k | X_i, D_i, Y_i; \theta) \stackrel{def}{=} \tilde{A}_{ki}(\theta) = \frac{L_i^C(\theta, a_k)}{\sum_{l=1}^K L_i^C(\theta, a_l)}, \quad (17)$$

for  $k = 1, \dots, K$ . In machine learning literature on classification (17) is called the “responsibility” of cluster  $k$  for unit  $i$ . We can use (17) to factor the  $i^{th}$  unit’s contribution to the complete data likelihood as

$$L_i^C(\theta; a_k) = \tilde{A}_{ki}(\theta) \left[ \sum_{l=1}^K L_i^C(\theta, a_l) \right].$$

This gives a re-arrangement of the lower bound (16) equal to

$$\begin{aligned} Q_i^*(\theta, q_i) &= \sum_{l=1}^K q_i(a_l) \ln \left( \frac{L_i^C(\theta, a_l)}{q_i(a_l)} \right) \\ &= \sum_{l=1}^K q_i(a_l) \ln \left( \frac{\tilde{A}_{li}(\theta) \left[ \sum_{m=1}^K L_i^C(\theta, a_m) \right]}{q_i(a_l)} \right) \\ &= -D_{\text{KL}}(q_i \| \tilde{A}_i) + \left[ \sum_{l=1}^K q_i(a_l) \right] \ln \left( \sum_{m=1}^K L_i^C(\theta, a_m) \right) \\ &= -D_{\text{KL}}(q_i \| \tilde{A}_i) + \ln \left( \sum_{m=1}^K L_i^C(\theta, a_m) \right), \end{aligned} \quad (18)$$

where  $D_{\text{KL}}(q_i \| \tilde{A}_i) = \sum_{l=1}^K q_i(a_l) \ln \left( \frac{q_i(a_l)}{\tilde{A}_{li}(\theta)} \right)$  is the Kullback-Leibler divergence of  $\tilde{A}_i$  from  $q_i$ .

Now consider, once again, the our maximization procedure:

1. Let  $\hat{\theta}^{(s)}$  for  $s = 0$  be an initial value for  $\theta$ .
2. **E-Step:** Set  $q(a_k) = \tilde{A}_{ki}(\hat{\theta}^{(s)})$  for  $k = 1, \dots, K$  and form the observed log-likelihood lower bound

$$\begin{aligned} Q(\theta, \hat{\theta}^{(s)}) &= \sum_{i=1}^N Q_i^*(\theta, \tilde{A}_i(\hat{\theta}^{(s)})) \\ &= \sum_{i=1}^N \left\{ \sum_{l=1}^K \tilde{A}_{li}(\hat{\theta}^{(s)}) \ln(L_i^C(\theta, a_l)) + \mathbf{S}(\tilde{A}_i(\hat{\theta}^{(s)})) \right\} \\ &= \sum_{i=1}^N \left\{ \mathbb{E} \left[ \ln(L_i^C(\theta, A)) \mid X_i, D_i, Y_i; \hat{\theta}^{(s)} \right] + \mathbf{S}(\tilde{A}_i(\hat{\theta}^{(s)})) \right\} \end{aligned} \quad (19)$$

where  $\mathbb{E} \left[ \ln(L_i^C(\theta, A)) \mid X_i, D_i, Y_i; \hat{\theta}^{(s)} \right]$  is the expected value of the  $i^{th}$  unit’s contribution to the complete data log-likelihood (given her observed encouragement, service

choice and earnings, and the current parameter value  $\hat{\theta}^{(s)}$ ) and  $\mathbf{S}(q) = -\sum_l q_l \ln q_l$  is the entropy of  $q$ .

3. **M-Step:** Choose  $\hat{\theta}^{(s+1)}$  to maximize  $Q(\theta, \hat{\theta}^{(s)})$  with respect to  $\theta$ . Note that since  $\mathbf{S}(\tilde{A}_i(\hat{\theta}^{(s)}))$  is constant in  $\theta$  this term is often omitted from the “Q-function” in practice (as was done above).
4. Repeat steps 2 and 3 until  $Q(\hat{\theta}^{(s+1)}, \hat{\theta}^{(s)}) \approx Q(\hat{\theta}^{(s)}, \hat{\theta}^{(s-1)})$  is small and/or until  $\hat{\theta}^{(s+1)} \approx \hat{\theta}^{(s)}$ .

Note that  $Q(\theta, \theta) = l_N^I(\theta)$ : after the E-Step the “Q-function” coincides with the observed log-likelihood (the Kullback-Leibler term is zero at  $q_i = \tilde{A}_i(\theta)$ ). We also have that the M-Step weakly increases the “Q-function”. Putting things together we have

$$l_N^I(\hat{\theta}^{(s+1)}) \geq Q(\hat{\theta}^{(s+1)}, \hat{\theta}^{(s)}) \geq Q(\hat{\theta}^{(s)}, \hat{\theta}^{(s)}) = l_N^I(\hat{\theta}^{(s)}). \quad (20)$$

From left-to-right the first inequality follows from (16), the second from the definition of maximization, and the third from (18) evaluated at  $q_i = \tilde{A}_i(\theta)$ . By virtue of (20) the observed log-likelihood  $Q(\hat{\theta}^{(s)}, \hat{\theta}^{(s)}) = l_N^I(\hat{\theta}^{(s)})$  is monotonically increasing in  $s$ . The EM algorithm will therefore find a local maximum (or saddle point) of the *observed log-likelihood* ((15)). Running the algorithm from a variety of starting points is advised.

## Further reading

The estimator introduced here is due to Imbens & Rubin (1997). While used frequently in applied statistics, its use in economics is less common, where method-of-moments estimators predominate (Angrist et al., 1996). One advantage of the model-based approach is the ease with which pre-treatment covariates are introduced into the analysis in a coherent way; something that is more difficult with the method-of-moments approach.

Murphy (2012, Chapter 12) provides a elementary introduction to the EM algorithm from a machine learning perspective. Gupta & Chen (2010) provide a survey with signal processing applications. Ruud (1991) provides a nice theoretical discussion with applications to discrete choice models common in econometrics. The EM algorithm has numerous applications in econometrics, particularly for discrete heterogeneity modeling.

## References

- Angrist, J. D. (1990). Lifetime earnings and the vietnam era draft lottery: evidence from social security administrative records. the. *American Economic Review*, 80(3), 313 – 336.
- Angrist, J. D., Imbens, G. W., & Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association*, 91(434), 444 – 455.
- Gupta, M. R. & Chen, Y. (2010). Theory and use of the em algorithm. *Foundations and Trends in Signal Processing*, 4(3), 223 – 296.
- Imbens, G. W. & Rubin, D. B. (1997). Estimating outcome distributions for compliers in instrumental variable models. *Review of Economic Studies*, 64(4), 555 – 574.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. Cambridge, MA: The MIT Press.
- Ruud, P. A. (1991). Extensions of estimation methods using the em algorithm. *Journal of Econometrics*, 49(3), 305 – 341.