

TA SESSION 6: HECKMAN SELECTION AND MULTINOMIAL PROBIT

Shuowen Chen¹

EC708: PhD Econometrics I (Spring 2020)

¹Parts of the materials are borrowed from Newey and McFadden (1994), Vella (1998), Puhani (2000), textbooks by Miranda and Fackler (2002), Train (2009) and Wooldridge (2010), and lecture slides by Gianluca Violante. I'm grateful to Yinuo Zhang for kindly sharing Violante's slides. I thank Iván Fernández-Val, Jean-Jacques Forneron and Hiroaki Kaido for helpful feedbacks.

OUTLINE

- ▶ Sample Selection
 - ▶ Motivations and Clarifications
 - ▶ Heckman Selection Model
 - ▶ Two-Step Estimators as GMM
- ▶ Probit Estimation
 - ▶ Comparison with Logit
 - ▶ Newton-Cotes method
 - ▶ Gaussian Quadrature
 - ▶ Accept-Reject Algorithm
 - ▶ GHK (Geweke-Hajivassiliou-Keane) Estimator
 - ▶ Importance Sampling Interpretation of GHK Estimator

RETURNS OF SCHOOLING ON WAGE RATE

- ▶ Research question: How does years of schooling affect wages?
- ▶ Why OLS of wages on yrs of schooling might be problematic? For some individuals the decision to work is not random
 1. Sample (working people) is unrepresentative of population of interest (all people who received schooling)²
 2. Some component of work decision relevant to wage determination
 3. If component is fully controlled by observable characteristics, can still run OLS
 4. Component might be unobservable: e.g., people who don't work were offered rates below their reservation wages
 5. Reservation wage is arguably correlated with ability, which determines the wage
 6. Failure to account for this correlation brings in endogeneity problem and leads to incorrect estimation

²Unless assume sample of working people is chosen randomly from the population

WHEN DOES SAMPLE SELECTION MATTER?

Sample Selection Bias arises whenever one examines a subsample and the **unobservable factors** determining inclusion in the subsample are **correlated with** the **unobservables** influencing the variable of primary interest (Vella, 1998)

HECKMAN SELECTION MODEL

Let N and n denote whole sample and subsample with observed dependent variable

$$y_i^* = x_i' \beta + \varepsilon_i \quad i = 1, \dots, N \quad (1)$$

$$d_i^* = z_i' \gamma + \nu_i \quad i = 1, \dots, N \quad (2)$$

$$d_i = \begin{cases} 1 & \text{if } d_i^* > 0 \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

$$y_i = y_i^* \times d_i \quad (4)$$

- ▶ Eq (1) is of primary interest
- ▶ Eq (2): reduced form for latent variable capturing sample selection
- ▶ Eq (3): whether the dependent variable is observed
- ▶ Eq (4): observed outcomes (data)
- ▶ (ε, ν) independent of z with zero mean
- ▶ ε and ν are correlated

MLE FOR HECKMAN SELECTION

Assumption: ε_i and ν_i are i.i.d distributed $\mathcal{N}(0, \Sigma)$, where

$$\Sigma = \begin{pmatrix} \sigma_{\varepsilon}^2 & \sigma_{\varepsilon\nu} \\ \sigma_{\varepsilon\nu} & \sigma_{\nu}^2 \end{pmatrix}$$

and (ε_i, ν_i) are independent of z_i .

Average log likelihood function:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \left\{ d_i \times \ln \left[\int_{-z_i'\gamma}^{\infty} \phi_{\varepsilon\nu}(y_i - x_i'\beta, \nu) d\nu \right] \right. \\ \left. + (1 - d_i) \times \left[\ln \int_{-z_i'\gamma}^{\infty} \int_{-\infty}^{\infty} \phi_{\varepsilon\nu}(\varepsilon, \nu) d\varepsilon d\nu \right] \right\}$$

- ▶ $\phi_{\varepsilon\nu}$: pdf of bivariate normal distribution
- ▶ Also known as Tobit type two model
- ▶ Fully efficient, but subject to misspecification
- ▶ Heckman (1976) proposes a simpler way for estimation

CORRECTION TERM OF HECKMAN SELECTION

Consider the version that assumes **bivariate normality**³. Note:

$$\mathbb{E}[y_i | z_i, d_i = 1] = x_i' \beta + \mathbb{E}[\varepsilon_i | z_i, d_i = 1], \quad i = 1, \dots, n$$

- ▶ $\mathbb{E}[\varepsilon_i | z_i, d_i = 1] = \frac{\sigma_{\varepsilon v}}{\sigma_v^2} \left\{ \frac{\phi(z_i' \gamma)}{\Phi(z_i' \gamma)} \right\}$
- ▶ Blue term: **inverse Mills ratio** (λ_i)

Remarks:

- ▶ If $\sigma_{\varepsilon v} = 0$, no correction term, sample selection is not a problem
- ▶ Estimate γ by Probit over the entire sample N by MLE
- ▶ Plug in $\hat{\gamma}$ for estimation of β (**two-step estimation**)
- ▶ Correction term plays the role of **control function**

³Other versions relax distribution assumptions or incorporate variable selections.

TWO-STEP ESTIMATION PROCEDURES & REMARKS

1. Estimate γ using the entire sample via MLE, compute $\hat{\lambda}$
2. Use the subsample with observed dependent variable to run OLS on the following specification:

$$y_i = x_i' \beta + \mu \hat{\lambda}_i + \eta_i,$$

where $\mu := \sigma_{\varepsilon\gamma} / \sigma_{\gamma}^2$.

Remarks:

- ▶ Need to adjust for standard errors (Greene, 1981)
- ▶ Can formulate as GMM and get standard errors easily
- ▶ **Identification concern:** the function mapping the single index $z_i' \gamma$ into inverse Mills ratio is **linear** for certain ranges of the index
 - ▶ Needs **exclusion restriction:** at least one variable in z_i is not in x_i
 - ▶ Otherwise identification of β relies on the **nonlinear part** of the inverse Mills ratio
 - ▶ Leads to weak identification and inflated second-step SE

TWO-STEP ESTIMATOR AS GMM

- ▶ Depends on some preliminary “first-step” estimator
- ▶ First-step estimation affects asymptotic variance of the second if consistency of the former affects that of the latter

Consider the following sample moments

$$\frac{1}{N} \sum_{i=1}^N g(w_i, \theta, \hat{\gamma}) = 0,$$

where θ is parameter of interest and $\hat{\gamma}$ is an estimator of γ , which satisfies the following sample moment condition

$$\frac{1}{N} \sum_{i=1}^N m(z_i, \gamma) = 0$$

Form $\tilde{g}(w, \theta, \gamma) = [m(z, \theta)', g(w, \theta, \gamma)']'$ and consider

$$\frac{1}{N} \sum_{i=1}^N \tilde{g}(w_i, \hat{\theta}, \hat{\gamma}) = 0$$

HECKMAN SELECTION MODEL AS TWO-STEP ESTIMATOR

Denote $\theta = (\beta', \mu)'$,

$$g(w, \theta, \gamma) = d \begin{bmatrix} x \\ \lambda(z' \gamma) \end{bmatrix} (y - x' \beta - \mu \lambda(z' \gamma)), \quad (5)$$

$$m(z, \gamma) = \lambda(z' \gamma) \Phi^{-1}(-z' \gamma) z (d - \Phi(z' \gamma)) \quad (6)$$

FOCs for OLS on the selected sample and probit estimation

Theorem 6.1 in Newey and McFadden (1994), can show that $\hat{\theta}$ and $\hat{\gamma}$ are asymptotically normal and $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V)$, where

$$V = G_{\theta}^{-1} \mathbb{E}[\{g(w, \theta_0, \gamma_0) + G_{\gamma} \psi(z)\} \{g(w, \theta_0, \gamma_0) + G_{\gamma} \psi(z)\}'] G_{\theta}^{-1'}$$

- ▶ $G_{\theta} = \mathbb{E}[\nabla_{\theta} g(w, \theta_0, \gamma_0)]$
- ▶ $G_{\gamma} = \mathbb{E}[\nabla_{\gamma} g(w, \theta_0, \gamma_0)]$
- ▶ $\psi(z) = -\mathbb{E}[\nabla_{\gamma} m(z, \gamma_0)]^{-1} m(z, \gamma_0)$

HOW FIRST-STEP ESTIMATION AFFECTS SECOND-STEP SE

- ▶ Suppose we only work with the following moment condition:

$$\frac{1}{N} \sum_{i=1}^N g(w_i, \theta, \gamma_0) = 0,$$

Then we have $Avar(\hat{\theta}) = G_{\theta}^{-1} \mathbb{E}[g(w, \theta_0, \gamma_0)g(w, \theta_0, \gamma_0)'] G_{\theta}^{-1}$,

- ▶ Simplified because we don't use first step $\hat{\gamma}$
- ▶ **Implication:** unless $G_{\gamma} = 0$, first-step estimation $\hat{\gamma}$ affects standard errors of $\hat{\theta}$
- ▶ When does $G_{\gamma} \neq 0$? **If inconsistency of $\hat{\gamma}$ leads to inconsistency of $\hat{\theta}$** (Newey & McFadden, 1994, Theorem 6.2)

FIRST-STEP ESTIMATION IN HECKMAN SELECTION

Denote $\lambda_v(v) = \frac{d\lambda(v)}{dv}$, G_γ in Heckman selection model is

$$G_\gamma = -\mu \mathbb{E} \left[d \left(\begin{matrix} x \\ \lambda(z' \gamma_0) \end{matrix} \right) \lambda_v(z' \gamma_0) z' \right],$$

- ▶ Generally nonzero unless $\mu = 0$, which implies $\sigma_{\varepsilon v} = 0$
- ▶ When does this happen? Sample selection doesn't matter

MORE ON ASYMPTOTIC VARIANCE

The correct asymptotic variance can be either larger or smaller than the one that ignores the first-step estimation

- ▶ Condition that correct asymptotic variance is larger:

$$\mathbb{E}[g(w, \theta_0, \gamma_0)m(z, \gamma_0)'] = 0$$

In this case $\mathbb{E}[g(w, \theta_0, \gamma_0)\psi(z)'] = 0$, and the correct variance is

$$G_\theta^{-1}\mathbb{E}[g(w, \theta_0, \gamma_0)g(w, \theta_0, \gamma_0)']G_\theta^{-1'} + G_\theta^{-1}G_\gamma\mathbb{E}[\psi(z)\psi(z)']G_\gamma'G_\theta^{-1'}$$

Heckman selection satisfies this condition since

$$\mathbb{E}[y - x'\beta_0 - \mu_0\lambda(z'\gamma_0)|x, d = 1, z] = 0$$

- ▶ Condition that correct asymptotic variance is smaller:

$$m(z, \gamma_0) = \nabla_\gamma \ln f(z|\theta_0, \gamma_0),$$

where f is likelihood of z . Somewhat rare to satisfy in practice.

CONSISTENT ASYMPTOTIC VARIANCE ESTIMATION

- ▶ The advantage of writing two-step estimation as a GMM is to get consistent asymptotic variance when $G_\gamma \neq 0$
- ▶ Recall asymptotic variance of $\hat{\theta}$:

$$V = G_\theta^{-1} \mathbb{E}[\{g(w, \theta_0, \gamma_0) + G_\gamma \psi(z)\} \{g(w, \theta_0, \gamma_0) + G_\gamma \psi(z)\}'] G_\theta^{-1'}$$

- ▶ Plugged-in approach:

$$\widehat{G}_\theta = \frac{1}{N} \sum_{i=1}^N \nabla_\theta g(w_i, \widehat{\theta}, \widehat{\gamma}), \quad \widehat{G}_\gamma = \frac{1}{N} \sum_{i=1}^N \nabla_\gamma g(w_i, \widehat{\theta}, \widehat{\gamma})$$

$$\widehat{g}_i = g(w_i, \widehat{\theta}, \widehat{\gamma}), \quad \widehat{m}_i = m(z_i, \widehat{\gamma}) \quad \widehat{\psi}_i = - \left[\frac{1}{N} \sum_{i=1}^N \nabla_\gamma m(z_i, \widehat{\gamma}) \right]^{-1} \widehat{m}_i$$

CONSISTENT ASYMPTOTIC VARIANCE ESTIMATION

The estimator of asymptotic variance of $\widehat{\theta}$ is

$$\widehat{V} = \widehat{G}_{\theta}^{-1} \left[\frac{1}{N} \sum_{i=1}^N (\widehat{g}_i + \widehat{G}_{\gamma} \widehat{\psi}_i) (\widehat{g}_i + \widehat{G}_{\gamma} \widehat{\psi}_i)' \right] \widehat{G}_{\theta}^{-1}$$

If moment functions are uncorrelated (recall this means **first-step estimation increases second-step variance**):

$$\mathbb{E}[g(w, \theta_0, \gamma_0) m(z, \gamma_0)'] = 0$$

Denote $\widehat{V}_{\gamma} = \frac{1}{N} \sum_{i=1}^N \widehat{\psi}_i \widehat{\psi}_i'$,

$$\widehat{V} = \widehat{G}_{\theta}^{-1} \left[\frac{1}{N} \sum_{i=1}^N \widehat{g}_i \widehat{g}_i' \right] \widehat{G}_{\theta}^{-1} + \widehat{G}_{\theta}^{-1} \widehat{G}_{\gamma} \widehat{V}_{\gamma} \widehat{G}_{\gamma}' \widehat{G}_{\theta}^{-1}$$

Can be applied to obtain variance estimator of Heckman selection

► Form

COMPARISON BETWEEN LOGIT AND PROBIT

- ▶ A standard logistic distribution has PDF:

$$f(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2}$$

and variance $\pi^2/3$. Looks very similar to normal distribution.

- ▶ (Multinomial) logit is widely used in empirical IO due to analytical forms and fast computations, but also limited:
 1. Can't represent random taste variation
 2. Restrictive substitution patterns due to IIA⁴
 3. Can't be used in panel if unobserved factors correlated over time for each individual
- ▶ Probit probabilities don't have closed-form and requires numerical approximations.
- ▶ We explain why this is the case and cover some methods.

▶ Graph

⁴Independence of irrelevant alternatives

MULTINOMIAL PROBIT: DISCRETE CHOICE UTILITY

Consider individual n 's additive utility of choosing one product among J options

$$U_{n,j} = V_{n,j} + \varepsilon_{n,j}$$

Assume $\varepsilon_n = (\varepsilon_{n,1}, \dots, \varepsilon_{n,J})$ is normally distributed with mean vector of zeros and var-cov matrix Ω . The probability of choose i is

$$\begin{aligned} P_{n,i} &= \Pr(V_{n,i} + \varepsilon_{n,i} > V_{n,j} + \varepsilon_{n,j} \forall j \neq i) \\ &= \int \mathbb{1}(V_{n,i} + \varepsilon_{n,i} > V_{n,j} + \varepsilon_{n,j} \forall j \neq i) \phi(\varepsilon_n) d\varepsilon_n, \end{aligned}$$

where $\phi(\varepsilon_n) = \frac{1}{(2\pi)^J |\Omega|^{1/2}} \exp\left(-\frac{1}{2} \varepsilon_n' \Omega^{-1} \varepsilon_n\right)$

The integration is over a vector, evaluated numerically. Essentially a numerical integration problem. Let's look at nonsimulation and simulation-based methods.


NEWTON-COTES METHOD: TRAPEZOID RULE

Consider the following integration problem: $\int_a^b f(x)dx$. Approximate f with a **piecewise linear polynomial** \tilde{f} whose integral is easy to compute:

$$\int_a^b f(x)dx \approx \int_a^b \tilde{f}(x)dx$$

1. Partition $[a, b]$ into n subintervals of equal length $h = \frac{b-a}{n}$ and endpoint nodes $x_k = a + kh$
2. For each node k , compute $y_k = f(x_k)$
3. Form a piecewise linear approximation of the function between successive points (x_k, x_{k+1}) :

$$f(x) \approx f(x_k) + \frac{f(x_{k+1}) - f(x_k)}{x_{k+1} - x_k}(x - x_k)$$

Simple and robust, increasing nodes from N to $M \times N$ reduces error by factors of M^2 

NEWTON–COTES METHOD: SIMPSON RULE

Instead of using piecewise linear polynomials, use **quadratic approximation**

- ▶ Form a piecewise quadratic approximation \tilde{f} that interpolates f at successive triplets (x_{k-1}, x_k, x_{k+1}) with quadratic functions
- ▶ Works better than trapezoid rule if f is smooth
- ▶ Works worse if f is non-differentiable at some points

▶ Graph

GAUSSIAN QUADRATURE

Consider a slightly more general integration problem: $\int_a^b f(x)w(x)dx$, where $w(\cdot)$ is a weight function (like *pdf*)

- ▶ Choose n nodes x_1, \dots, x_n and weights⁵ w_1, \dots, w_n to match the $2n$ conditions: $\int_a^b x^k w(x)dx = \sum_{i=1}^n w_i x_i^k, \quad k = 0, \dots, 2n-1$
- ▶ Integral approximation is thus $\int_a^b f(x)w(x)dx \approx \sum_{i=1}^n w_i f(x_i)$
- ▶ When $w(x) = 1$, called Gauss–Legendre quadrature
- ▶ When $w(x)$ is pdf, $\sum_{i=1}^n w_i x_i^k = EX^k$: discretize continuous r.v. with mass x_i and prob w_i and match moments
- ▶ For smooth integrands, converges exponentially fast as $n \uparrow$
- ▶ For more details, refer to *Numerical Recipes* by Press et al. (2007)

⁵Some common weight functions: uniform, normal, gamma, beta, etc.

ACCEPT-REJECT ALGORITHM

$$P_{n,i} = \int \mathbb{1}(V_{n,i} + \varepsilon_{n,i} > V_{n,j} + \varepsilon_{n,j} \forall j \neq i) \phi(\varepsilon_n) d\varepsilon_n$$

1. Draw $\varepsilon_n = (\varepsilon_{n,1}, \dots, \varepsilon_{n,j})$ from a normal density with zero mean and var-cov Ω
2. Calculate $U_{n,j} = V_{n,j} + \varepsilon_{n,j} \forall j$
3. Determine if $U_{n,i} > U_{n,j} \forall j \neq i$. If yes, count as **accept**, otherwise count as **reject**
4. Repeat steps 1–3 many times (S)
5. Simulated probability is the proportion of draws that are **accepted**

Remarks:

- ▶ How to draw from a joint normal density? ▶ Procedure
- ▶ Simulated prob might be zero for any finite number of draws, especially if true prob is low
- ▶ Simulated prob not smooth, hinders numerical optimization

SMOOTHED ACCEPT-REJECT SIMULATORS

Accept-Reject algorithm is not smooth because it uses 0-1 indicator, replace with a **smooth and strictly positive** function. Following McFadden (1989), we use logit function.

1. Draw $\varepsilon_n = (\varepsilon_{n,1}, \dots, \varepsilon_{n,j})$ from a normal density with zero mean and var-cov Ω
2. Calculate $U_{n,j} = V_{n,j} + \varepsilon_{n,j} \forall j$
3. Calculate $M = \frac{\exp(U_{n,i})/\lambda}{\sum_j \exp(U_{n,j})/\lambda}$
4. Repeat steps 1-3 many times (S)
5. Simulated probability is the average of M_s : $\widehat{P}_{n,i} = \frac{1}{S} \sum_{s=1}^S M_s$

Remark:

- ▶ λ : degree of smoothing, approaches Accept-Reject as $\lambda \rightarrow 0$

GEWEKE–HAJIVASSILIOU–KEANE (GHK) SIMULATOR

- ▶ Based on sampling from recursive truncated normals after a Cholesky transformation ▶ Cholesky Decomposition
- ▶ First suggested by Geweke (1989), independently developed by Hajivassiliou (1992) and Keane (1994)
- ▶ Operates on utility differences: if we want to simulate $P_{n,i}$, need to subtract $U_{n,i}$ from the other utilities

Consider a three-alternative case and simulate $P_{n,1}$:

$$\widetilde{U}_{n,j,1} \equiv U_{n,j} - U_{n,1} = (V_{n,j} - V_{n,1}) + (\varepsilon_{n,j} - \varepsilon_{n,1}) \equiv \widetilde{V}_{n,j,1} + \widetilde{\varepsilon}_{n,j,1}$$

$\varepsilon_n = (\varepsilon_{n,1}, \varepsilon_{n,2}, \varepsilon_{n,3})'$ has distribution $\mathcal{N}(0, \Omega)$; $\widetilde{\varepsilon}_{n,1} = (\widetilde{\varepsilon}_{n,2,1}, \widetilde{\varepsilon}_{n,3,1})'$ has distribution $\mathcal{N}(0, \widetilde{\Omega}_1)$: $\widetilde{\Omega}_1 = A_1 \Omega A_1'$, where

$$A_1 = \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix}$$

GHK SIMULATOR: PREPARATION

- ▶ Consider the lower-triangular Cholesky factor of $\tilde{\Omega}_1$,

$$L = \begin{pmatrix} c_{aa} & 0 \\ c_{ab} & c_{bb} \end{pmatrix}$$

- ▶ Rewrite $(\tilde{\varepsilon}_{n,2,1}, \tilde{\varepsilon}_{n,3,1})$ as

$$\tilde{\varepsilon}_{n,2,1} = c_{aa}\eta_1, \quad \tilde{\varepsilon}_{n,3,1} = c_{ab}\eta_1 + c_{bb}\eta_2$$

where η_1 and η_2 are i.i.d. standard normal

- ▶ Hence we have

$$\tilde{U}_{n,2,1} = \tilde{V}_{n,2,1} + c_{aa}\eta_1, \quad \tilde{U}_{n,3,1} = \tilde{V}_{n,3,1} + c_{ab}\eta_1 + c_{bb}\eta_2$$

- ▶ Prob that option 1 is chosen:

$$P_{n,1} = \Pr(\tilde{U}_{n,2,1} < 0 \ \& \ \tilde{U}_{n,3,1} < 0)$$

- ▶ **Why need Cholesky decomposition?** $(\tilde{\varepsilon}_{n,2,1}, \tilde{\varepsilon}_{n,3,1})$ are correlated, hard to evaluate numerically, but η_1 and η_2 are easy to draw

GHK SIMULATOR: FORM TO WORK WITH

Using Bayes rule,

$$\begin{aligned}P_{n,1} &= Pr(\tilde{V}_{n,2,1} + c_{aa}\eta_1 < 0 \ \& \ \tilde{V}_{n,3,1} + c_{ab}\eta_1 + c_{bb}\eta_2 < 0) \\&= Pr(\tilde{V}_{n,2,1} + c_{aa}\eta_1 < 0)Pr(\tilde{V}_{n,3,1} + c_{ab}\eta_1 + c_{bb}\eta_2 < 0 | \tilde{V}_{n,2,1} + c_{aa}\eta_1 < 0) \\&= Pr\left(\eta_1 < -\frac{\tilde{V}_{n,2,1}}{c_{aa}}\right)Pr\left(\eta_2 < -\frac{\tilde{V}_{n,3,1} + c_{ab}\eta_1}{c_{bb}} \middle| \eta_1 < -\frac{\tilde{V}_{n,2,1}}{c_{aa}}\right) \\&= \Phi\left(-\frac{\tilde{V}_{n,2,1}}{c_{aa}}\right) \int_{-\infty}^{-\frac{\tilde{V}_{n,2,1}}{c_{aa}}} \Phi\left(-\frac{\tilde{V}_{n,3,1} + c_{ab}\eta_1}{c_{bb}}\right) \bar{\phi}(\eta_1) d\eta_1\end{aligned}$$

where

$$\bar{\phi}(\eta_1) = \begin{cases} \frac{\phi(\eta_1)}{\Phi(-\tilde{V}_{n,2,1}/c_{aa})} & -\infty < \eta_1 < -\tilde{V}_{n,2,1}/c_{aa} \\ 0 & o.w \end{cases}$$

- ▶ Blue term is easy to compute
- ▶ Red term: **GHK punchline** (approximate integral by simulations)

GHK SIMULATOR: ALGORITHM

Procedures to approximate the red term:

1. Draw η_1 from truncated normal density $\bar{\phi}$
2. Compute $\Phi\left(-\frac{\tilde{V}_{n,3,1}+c_{ab}\eta_1}{c_{bb}}\right)$
3. Repeat 1–2 many times and take average

How to draw from a truncated univariate distribution? 

GHK Algorithm:

1. Compute $\Phi\left(-\frac{\tilde{V}_{n,2,1}}{c_{aa}}\right)$
2. Conduct procedures to approximate the red term, denote the outcome as B
3. Simulated probability $\hat{P}_{n,1} = \Phi\left(-\frac{\tilde{V}_{n,2,1}}{c_{aa}}\right)B$

Remarks:

- ▶ Applicable to simulate other choice probabilities, but with different $\tilde{\Omega}_i$ and hence Cholesky factors
- ▶ Can be generalized to many-alternative problems

DETOUR: IMPORTANCE SAMPLING

Suppose r.v. u has a density $f(u)$ that's hard to draw from directly, but there is another density $g(u)$ that is easy to draw from. We can obtain draws from f using the following procedure:

1. Draw from g and denote it as $u^{(1)}$
2. Weight the draw by $f(u^{(1)})/g(u^{(1)})$
3. Repeat 1–2 many times
4. The set of weighted draws is equal to a set of draws from f

CDF of weighted draws from g equals that from f :

$$\int \frac{f(u)}{g(u)} \mathbb{1}_{\{u < m\}} g(u) du = \int_{-\infty}^m \frac{f(u)}{g(u)} g(u) du = F(m)$$

Requires two conditions:

- ▶ Support of g covers that of f : any possible u of f drawable from g
- ▶ $\forall u, \mathbb{E}\left(\frac{f(u)}{g(u)}\right) < \infty$

APPROXIMATE INTEGRAL VIA IMPORTANCE SAMPLING

Suppose we want to calculate the following:

$$\int t(\varepsilon)f(\varepsilon)d\varepsilon,$$

where $t(\varepsilon)$ is a function of ε and f is the density that is hard to draw from directly. Suppose an alternative density g is easy to draw from. Rewrite integral to be

$$\int \left[t(\varepsilon)\frac{f(\varepsilon)}{g(\varepsilon)} \right] g(\varepsilon)d\varepsilon$$

Draw ε from g and evaluate $t\frac{f}{g}$, repeat S times and take average

Remarks:

- ▶ In practice g should have a larger tail than f
- ▶ Compute the effective sample size $\frac{1}{\sum_i^S w_i^2}$, where w_i is the normalized importance weights, which can be very low (as in particle filters)

ANOTHER LOOK OF GHK

Recall the choice 1 probability:

$$P_{n,1} = \int \mathbb{1}(\eta \in \mathcal{B}) f(\eta) d\eta,$$

where $\mathcal{B} = \{\eta : \tilde{U}_{n,2,1} < 0 \text{ \& } \tilde{U}_{n,3,1} < 0\}$ and $f(\eta) = \phi(\eta_1)\phi(\eta_2)$

- ▶ Accept-Reject algorithm **draws from f** and evaluate if $\mathbb{1}(\eta \in \mathcal{B})$ is accepted given the draws
- ▶ GHK **doesn't draw from f** but from **g , truncated normal**

$$g(\eta) = \begin{cases} \frac{\phi(\eta_1)}{\Phi(-\tilde{V}_{n,2,1})/c_{aa}} \times \frac{\phi(\eta_2)}{\Phi(-(\tilde{V}_{n,3,1} + c_{ab}\eta_1))/c_{bb}} & \eta \in \mathcal{B} \\ 0 & \eta \notin \mathcal{B} \end{cases}$$

- ▶ GHK only draws from densities that are consistent with player choosing option 1

ANOTHER LOOK OF GHK

Recall in the GHK algorithm, for each draw η we essentially calculate the following:

$$\widehat{P}_{n,1}(\eta) = \Phi\left(-\frac{\widetilde{V}_{n,2,1}}{c_{aa}}\right) \times \Phi\left(-\frac{\widetilde{V}_{n,3,1} + c_{ab}\eta_1}{c_{bb}}\right)$$

which is the denominator of g for $\eta \in \mathcal{B}$, so rewrite g as

$$g(\eta) = \begin{cases} \frac{f(\eta)}{\widehat{P}_{n,1}(\eta)} & \eta \in \mathcal{B} \\ 0 & \eta \notin \mathcal{B} \end{cases}$$

IMPORTANCE SAMPLING INTERPRETATION OF GHK

Hence we have

$$\begin{aligned}P_{n,1} &= \int \mathbb{1}(\eta \in \mathcal{B}) f(\eta) d\eta \\&= \int \mathbb{1}(\eta \in \mathcal{B}) \frac{f(\eta)}{g(\eta)} g(\eta) d\eta \\&= \int \mathbb{1}(\eta \in \mathcal{B}) \frac{f(\eta)}{\frac{f(\eta)}{\widehat{P}_{n,1}(\eta)}} g(\eta) d\eta \\&= \int \mathbb{1}(\eta \in \mathcal{B}) \widehat{P}_{n,1}(\eta) g(\eta) d\eta \\&= \int \widehat{P}_{n,1}(\eta) g(\eta) d\eta\end{aligned}$$

Remarks:

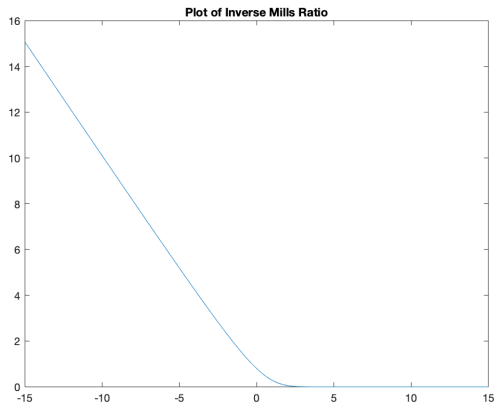
- ▶ Last equality is due to $g(\eta) > 0$ only when $\mathbb{1}(\eta \in \mathcal{B}) = 1$
- ▶ $\widehat{P}_{n,1}(\eta)$ is the weight imposed on draws from g
- ▶ GHK replaces $\mathbb{1}(\eta \in \mathcal{B})$ with smoothed $\widehat{P}_{n,1}(\eta)$

VARIANCE OF ERRORS IN THE SECOND STEP

$$\text{Var}(\eta_i) = \sigma_{\varepsilon}^2 - \frac{\sigma_{\varepsilon\nu}^2}{\sigma_{\nu}^2} \left[\frac{z_i'\gamma}{\sigma_{\nu}} \cdot \frac{\phi(z_i'\gamma)}{\Phi(z_i'\gamma)} + \left(\frac{\phi(z_i'\gamma)}{\Phi(z_i'\gamma)} \right)^2 \right]$$

◀ Back

ILLUSTRATION OF INVERSE MILLS RATIO



HECKMAN SELECTION VARIANCE ESTIMATOR

Denote $W_i = d_i[x_i', \lambda(z_i'\gamma_0)]'$ and $\widehat{W}_i = d_i[x_i', \lambda(z_i'\widehat{\gamma})]'$, then

$$G_\theta = -\mathbb{E}[d_i W_i W_i'], \quad G_\gamma = -\mu_0 \mathbb{E}[d_i \lambda_\nu(z_i'\gamma_0) W_i z_i'],$$

with sample analog

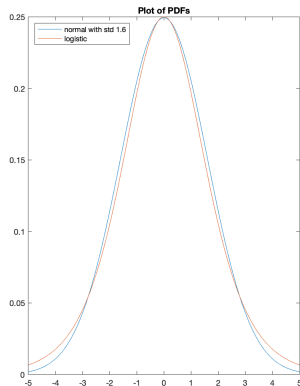
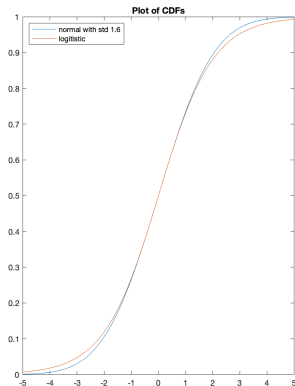
$$\widehat{G}_\theta = -\frac{1}{N} \sum_{i=1}^N \widehat{W}_i \widehat{W}_i', \quad \widehat{G}_\gamma = -\frac{1}{N} \widehat{\mu} \sum_{i=1}^N \lambda_\nu(z_i'\widehat{\gamma}) \widehat{W}_i z_i'$$

Denote $\widehat{\eta}_i = y_i - \widehat{W}_i'(\widehat{\beta}', \widehat{\mu})'$, the variance estimator is thus

$$\widehat{V} = N \left(\sum_{i=1}^N \widehat{W}_i \widehat{W}_i' \right)^{-1} \sum_{i=1}^N \widehat{W}_i \widehat{W}_i \widehat{\eta}_i^2 \left(\sum_{i=1}^N \widehat{W}_i \widehat{W}_i' \right)^{-1} + \widehat{\Pi} \widehat{V}_\gamma \widehat{\Pi}'$$

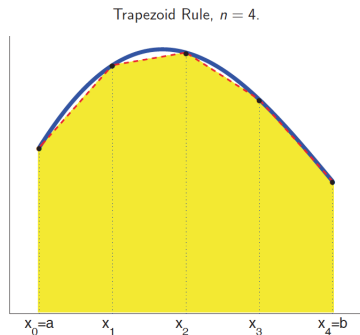
- ▶ Blue: sum of White estimator for least squares
- ▶ Red: correction term for first step. \widehat{V}_γ : *variance estimator of first step probit*; $\widehat{\Pi} = \widehat{G}_\theta^{-1} \widehat{G}_\gamma$.

PLOTS OF LOGIT AND NORMAL DISTRIBUTIONS



- I picked $\sigma = 1.6$ for normality to fit the two distributions.

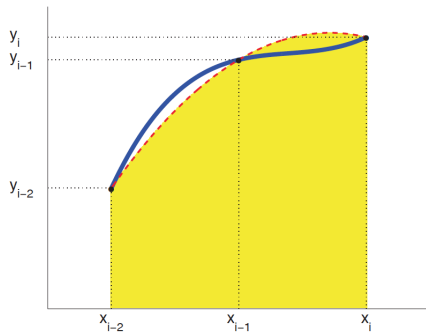
ILLUSTRATION OF NEWTON–COTES TRAPEZOID



Why trapezoid? The area under the piecewise linear approximation for subinterval k is

$$\int_{x_k}^{x_{k+1}} \tilde{f}(x) dx \approx \left[\frac{f(x_{k+1}) + f(x_k)}{2} \right] h$$

ILLUSTRATION OF NEWTON-COTES SIMPSON



◀ Back

CHOLESKY DECOMPOSITION

Consider a symmetric and positive-definite square matrix A . We can construct a lower triangular matrix L whose transpose serves as the upper triangular part:

$$L \cdot L' = A$$

Remark: Cholesky decomposition is an efficient way to check if a matrix is truly positive-definite. An alternative way is to check if the minimum eigenvalue is positive. [◀ Back](#)

DRAWING FROM JOINT NORMAL DENSITY

Consider a J -vector joint normal distribution $\mathcal{N}(b, \Omega)$. The Cholesky decomposition is

$$L \cdot L' = \Omega$$

1. Take J draws from a standard normal and denote as $v = (v_1, \dots, v_J)'$

2. Calculate $\varepsilon = b + Lv$

ε is normally distributed, has mean b and covariance Ω [◀ Back](#)

DRAWING FROM TRUNCATED UNIVARIATE DENSITY

Consider a r.v ranging from a to b with density proportional to $f(\varepsilon)$ within the range. In other words:

$$k = \int_a^b f(\varepsilon) d\varepsilon = F(b) - F(a)$$

The r.v has density $\frac{1}{k}f(\varepsilon)$ for $a \leq \varepsilon \leq b$ and 0 otherwise. The following procedure draws a value from this truncated density.

1. Draw μ from standard uniform $\mathcal{U}[0, 1]$
2. Calculate weighted average $\bar{\mu} = (1 - \mu)F(a) + \mu F(b)$
3. Calculate $\varepsilon = F^{-1}(\bar{\mu})$

The draw of μ determines how far to go between a and b 