

1. Linear Models

a. Regression

Think of this as the empirical analysis of the relationship between y (outcome variable) and \mathbf{x} (vector of regressors). \rightarrow convenient to think of regression as directed to the prediction of y given \mathbf{x} .

Define

$$\begin{aligned}\hat{y} &\equiv \textbf{predictor of } y \text{ as a function of } \mathbf{x} \\ \hat{u} &\equiv y - \hat{y} \\ &\equiv \textbf{prediction error} \\ L(\hat{u}) &= L(y - \hat{y}) \\ &\equiv \text{loss associated with } \hat{u} \\ &\quad (\text{differences in } \textbf{loss functions } L(\cdot) \text{ imply different regression models})\end{aligned}$$

Assume

$$\begin{aligned}L(0) &= 0 \\ L(u) &\geq 0, \quad \forall u \in \mathbb{R} \\ L(u) &\text{ is increasing in } |u|\end{aligned}$$

Assume (y, \hat{y}) is a random vector; goal is to minimize the (conditional) expected loss.

$$E[L(\hat{u})|\mathbf{x}]$$

In Economics, often no obvious choice of $L(\cdot)$, so usually specify quadratic loss, i.e.

$$L(u) = u^2 \Rightarrow E[\hat{u}^2|\mathbf{x}]$$

i.e. mean squared prediction error (MSPE).

In this case, the optimal predictor of y is the conditional mean $E[y|\mathbf{x}]$. One usually imposes a parametric model on $E[y|\mathbf{x}]$, i.e. one assumes that \exists a function $g(\mathbf{x}, \boldsymbol{\beta})$ and a vector $\boldsymbol{\beta}$ such that

$$P[E[y|\mathbf{x}] = g(\mathbf{x}, \boldsymbol{\beta})] = 1$$

Suppose we have observations

$$(\mathbf{x}_i^\top, y_i), \quad i = 1, \dots, n$$

Then the optimal predictor is

$$\begin{aligned}\hat{y} &= g(\mathbf{x}, \hat{\boldsymbol{\beta}}), \quad \text{where } \hat{\boldsymbol{\beta}} \text{ minimizes} \\ \sum_{i=1}^n L(\hat{u}_i) &= \sum_{i=1}^n \hat{u}_i^2 = \sum_{i=1}^n (y_i - g(\mathbf{x}_i, \boldsymbol{\beta}))^2\end{aligned}$$

i.e. the sum of squared residuals, i.e. $\hat{\boldsymbol{\beta}}$ is computed by nonlinear least squares (NLLS). If $g(\mathbf{x}, \boldsymbol{\beta})$ is actually linear in \mathbf{x} , $\boldsymbol{\beta}$, i.e. if $g(\mathbf{x}, \boldsymbol{\beta}) = \mathbf{x}^\top \boldsymbol{\beta}$, then the optimal predictor is

$$\hat{y} = \mathbf{x}^\top \hat{\boldsymbol{\beta}}$$

Where $\hat{\beta}$ is computed by ordinary least squares (OLS) (Legendre, 1805).

By contrast, under absolute error loss (i.e. $L(u)=|u|$), the optimal predictor is the conditional median of y given \mathbf{x} , i.e.

$$F_{y|\mathbf{x}}^{-1}(0.5)$$

If one additionally assumes linearity in the conditional median, i.e. if $F_{y|\mathbf{x}}^{-1}(0.5) = \mathbf{x}^\top \beta$, then the optimal predictor is $\hat{y} = \mathbf{x}^\top \hat{\beta}$, where $\hat{\beta}$ minimizes

$$\sum_{i=1}^n L(u_i) = \sum_{i=1}^n |y_i - \mathbf{x}_i^\top \beta|$$

i.e. the sum of absolute residuals, i.e. $\hat{\beta}$ is the least absolute deviations (LAD) estimator (Boscovich, 1757).

Go back to prediction under quadratic error loss. If $E[y|\mathbf{x}] = \mathbf{x}^\top \beta$ (w.p. 1), then β is of intrinsic interest. It has a causal interpretation (think a one unit increase in a component of \mathbf{x} implies a change in the optimal prediction of y equal to the corresponding component of β). So, OLS estimates an intrinsically important quantity when $E[y|\mathbf{x}]$ is linear.

If $E[y|\mathbf{x}]$ is nonlinear, might not be clear what OLS actually estimates.

Differentiate w.r.t. β the expected loss

$$E[(y - \mathbf{x}^\top \beta)^2]$$

and set the derivative to $\mathbf{0}$.

$$\begin{aligned} -2E[\mathbf{x}(y - \mathbf{x}^\top \beta)] &= \mathbf{0} \\ \Leftrightarrow \beta &= (E[\mathbf{x}\mathbf{x}^\top])^{-1}E[\mathbf{x}y] \end{aligned}$$

which is consistently estimated by OLS.

So, OLS estimates the best linear predictor of y under quadratic loss. Now suppose

$$\mathbf{x} = [1 \quad \tilde{\mathbf{x}}^\top]^\top, \quad \beta = [\beta_1 \quad \tilde{\beta}^\top]^\top$$

Then, $\mathbf{x}^\top \beta$ “has an intercept” since $\mathbf{x}^\top \beta = \beta_1 + \tilde{\mathbf{x}}^\top \tilde{\beta}$. Then the first-order conditions for minimization of expected loss can be written as

$$\begin{aligned} -2E[y - \beta_1 - \tilde{\mathbf{x}}^\top \tilde{\beta}] &= 0 \\ -2E[\tilde{\mathbf{x}}(y - \beta_1 - \tilde{\mathbf{x}}^\top \tilde{\beta})] &= \mathbf{0} \end{aligned}$$

This implies the existence of an unrestricted linear model for y since we can write

$$y = \beta_1 + \tilde{\mathbf{x}}^\top \tilde{\beta} + u, \quad \text{where } E[u] = 0 \text{ and } \text{Cov}[\tilde{\mathbf{x}}, u] = \mathbf{0}$$

Solve the above FOCs to get

$$\begin{aligned} \beta_1 &= E[y] - E[\tilde{\mathbf{x}}^\top] \tilde{\beta} \\ \tilde{\beta} &= (\text{Var}[\tilde{\mathbf{x}}])^{-1} \text{Cov}[\tilde{\mathbf{x}}, y] \end{aligned}$$

which are consistently estimable by OLS.

Also, for any combination of regressors \mathbf{x} and regressand y , one can (always) specify the linear projection model

$$y = \beta_1 + \tilde{\mathbf{x}}^\top \tilde{\beta} + u, \quad \text{where } E[u] = 0 \text{ and } \text{Cov}[\tilde{\mathbf{x}}, u] = \mathbf{0}$$

and where $\beta_1, \tilde{\beta}$ are given in the solutions of FOCs.

Note: Neither β_1 nor $\tilde{\beta}$ has a causal interpretation unless one imposes the additional restriction that $E[u|\mathbf{x}]=0$ (in which case we have a linear regression model).

b. OLS for linear regression

Observe

$$(\mathbf{x}_i^\top, y_i), \quad i = 1, \dots, n$$

where each \mathbf{x}_i is d-variate. Assume

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + u_i, \quad i = 1, \dots, n$$

where $\boldsymbol{\beta} \in \mathbb{R}^d$ is unknown and each u_i is random and unobserved. Write

$$\mathbf{y} \equiv \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} \equiv \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}, \quad \mathbf{u} \equiv \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}, \quad \boldsymbol{\beta} \equiv \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}$$

Then, $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$. The OLS estimator $\hat{\boldsymbol{\beta}}$ is defined as

$$\hat{\boldsymbol{\beta}} = \underset{\mathbf{b} \in \mathbb{R}^d}{\operatorname{argmax}} (\mathbf{y} - \mathbf{X}\mathbf{b})^\top (\mathbf{y} - \mathbf{X}\mathbf{b})$$

i.e. the minimizer of the sum of squared residuals. Differentiate the objective function w.r.t. \mathbf{b} , solve the FOCs to get

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

(this assumes $\mathbf{X}^\top \mathbf{X}$ is invertible; otherwise replace $(\mathbf{X}^\top \mathbf{X})^{-1}$ with a generalized inverse)

Proposition (asymptotic distribution of OLS) Assume

$$(A1) \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$$

$$(A2) E[\mathbf{u}|\mathbf{X}] = \mathbf{0}$$

$$E[\mathbf{u}\mathbf{u}^\top | \mathbf{X}] = \boldsymbol{\Omega}, \quad \text{where } \boldsymbol{\Omega} = \operatorname{diag}[\sigma_i^2] = \begin{bmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_n^2 \end{bmatrix} \text{ for constants } \sigma_1, \dots, \sigma_n > 0$$

$$(A3) \mathbf{X} \text{ has full column rank (so } \mathbf{X}^\top \mathbf{X} \text{ is nonsingular)}$$

$$(A4) \boldsymbol{\Sigma}_{\mathbf{xx}} \equiv \operatorname{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{X} \text{ exists and is finite and nonsingular}$$

$$(A5) \frac{1}{\sqrt{n}} \mathbf{X}^\top \mathbf{u} \xrightarrow{d} N_d(\mathbf{0}, \mathbf{S}), \quad \text{where } \mathbf{S} = \operatorname{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^\top \mathbf{u}\mathbf{u}^\top \mathbf{X}$$

Then, the OLS estimator $\hat{\boldsymbol{\beta}}$ is consistent and also satisfies

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &\xrightarrow{p} \boldsymbol{\beta} \\ \sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) &\xrightarrow{d} N_d(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1} \mathbf{S} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1}) \end{aligned}$$

Note

1.

$$\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top, \quad \text{so under A4,} \quad \boldsymbol{\Sigma}_{\mathbf{xx}} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[\mathbf{x}_i \mathbf{x}_i^\top]$$

2.

$$\frac{1}{n} \mathbf{X}^\top \mathbf{u} \mathbf{u}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n u_i^2 \mathbf{x}_i \mathbf{x}_i^\top, \quad \text{so under A5,} \quad \mathbf{S} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[u_i^2 \mathbf{x}_i \mathbf{x}_i^\top]$$

3. A1 implies that

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} = \boldsymbol{\beta} + \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1} \frac{1}{n} \mathbf{X}^\top \mathbf{u}$$

So consistency is straightforward via A4, A5 and the continuous mapping theorem.

4. A2 implies that

$$E[\mathbf{y}|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$$

which implies that $\boldsymbol{\beta}$ has a causal interpretation; A2 also permits (conditional) heteroscedasticity.

5. A3 rules out perfect multicollinearity, so $\mathbf{X}^\top \mathbf{X}$ is invertible.

6. A4, A5 jointly imply that $\hat{\boldsymbol{\beta}}$ is root-n-consistent, i.e. that $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ is stochastically bounded as $n \rightarrow \infty$, i.e. $O_p(1)$.

To see this; because we assume

$$\frac{1}{n} \mathbf{X}^\top \mathbf{X} \xrightarrow{p} \boldsymbol{\Sigma}_{\mathbf{xx}}$$

We have

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \sqrt{n} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1} \frac{1}{n} \mathbf{X}^\top \mathbf{u} = \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} \right)^{-1} \frac{1}{\sqrt{n}} \mathbf{X}^\top \mathbf{u}$$

which is asymptotically normal mean zero (i.e. not just stochastically bounded) by A4, A5 and Slutsky's theorem.

7. We've seen that $\hat{\boldsymbol{\beta}}$ has an approximate $N_d(\boldsymbol{\beta}, V[\hat{\boldsymbol{\beta}}])$ -distribution, where

$$V[\hat{\boldsymbol{\beta}}] = \frac{1}{n} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1} \mathbf{S} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1}, \quad \text{when } n \text{ is large}$$

Call $V[\hat{\boldsymbol{\beta}}]$ the asymptotic variance of $\hat{\boldsymbol{\beta}}$.

Observe $\{(\mathbf{x}_i^\top, y_i): i=1, \dots, n\}$.

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \text{and } \mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$$

Assume $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + u_i$ ($i=1, \dots, n$).

$$\text{OLS } \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

We gave assumptions that

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N_d(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1} \mathbf{S} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1})$$

i.e. $\hat{\boldsymbol{\beta}}$ has an approximate $N_d(\boldsymbol{\beta}, n^{-1} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1} \mathbf{S} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1})$ -distribution when n is large.

Call $\mathbf{V}[\hat{\boldsymbol{\beta}}] \equiv \frac{1}{n} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1} \mathbf{S} \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1}$ the asymptotic variance of $\hat{\boldsymbol{\beta}}$. The formulation of tests/confidence regions regarding $\boldsymbol{\beta}$ requires the (consistent estimation) of $\mathbf{V}[\hat{\boldsymbol{\beta}}]$.

Note that

- $\boldsymbol{\Sigma}_{\mathbf{xx}}$ is estimable via $\mathbf{S}_{\mathbf{xx}} \equiv n^{-1} \mathbf{X}^\top \mathbf{X}$ (in view of Assumption 4, hereafter A4).
- White (1980) proposed to estimate \mathbf{S} via $\hat{\mathbf{S}} \equiv \frac{1}{n} \mathbf{X}^\top \hat{\boldsymbol{\Omega}} \mathbf{X}$ where $\hat{\boldsymbol{\Omega}} = \text{diag}[\hat{u}_i^2: i = 1, \dots, n]$ and where $\hat{u}_i = y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$.
- Combining $\mathbf{S}_{\mathbf{xx}}$, $\hat{\mathbf{S}}$ appropriately gives an estimate of $\mathbf{V}[\hat{\boldsymbol{\beta}}]$.

$$\hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}] \equiv \frac{1}{n} \mathbf{S}_{\mathbf{xx}}^{-1} \hat{\mathbf{S}} \mathbf{S}_{\mathbf{xx}}^{-1} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\boldsymbol{\Omega}} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}$$

- If u_1, \dots, u_n are homoskedastic, so $\boldsymbol{\Omega} \equiv E[\mathbf{u}\mathbf{u}^\top | \mathbf{X}] = \sigma^2 \mathbf{I}_n$, then $\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} = \sigma^2 \mathbf{X}^\top \mathbf{X}$ and $\mathbf{S} = \sigma^2 \boldsymbol{\Sigma}_{\mathbf{xx}}$, which implies $\mathbf{V}[\hat{\boldsymbol{\beta}}] = \frac{1}{n} \sigma^2 \boldsymbol{\Sigma}_{\mathbf{xx}}^{-1}$. This leads to the default (or “homoskedasticity”).

$$\text{OLS Variance Estimator } \tilde{\mathbf{V}}[\hat{\boldsymbol{\beta}}] \equiv \frac{1}{n} s^2 \mathbf{S}_{\mathbf{xx}}^{-1} = s^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

$$\text{where } s^2 = \frac{1}{n-d} \sum_{i=1}^n \hat{u}_i^2$$

(i.e. the square of the regression standard error)

($\tilde{\mathbf{V}}[\hat{\boldsymbol{\beta}}]$ is inconsistent under heteroskedasticity)

How realistic are the assumptions in the previous proposition? Consider six “primitive” assumptions (White, 1980) that imply A1–A5 in the previous proposition.

Suppose $\mathbf{x}_i = [x_{i1} \ x_{i2} \ \dots \ x_{id}]^\top$ for each i .

(A1–A5 are used for the assumptions of Lecture 1. Though the professor used A1–A6 in the class, AA1–AA6 are used for the assumptions of Lecture 2 in order to distinguish them.)

AA1 The observations $[\mathbf{x}_i^\top \ y_i]$ ($i=1, \dots, n$) are independent.

- This allows for the data to be generated by stratified random sampling (i.e. split the population into strata, make random draws within each stratum. Typically see oversampling of certain strata).
- AA1 (obviously) rules out most time-series settings, as well as the “spatial” clustering of observations.

AA2 $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + u_i$ for each i .

- This assumes correctness of the linear model specification. Failure of AA2 implies that OLS’s consistent only for the best linear predictor.

AA3 Each vector \mathbf{x}_i of regressors is possibly random with $E[\mathbf{x}_i \mathbf{x}_i^\top]$ finite. In addition, $E[|\mathbf{x}_{ij} \mathbf{x}_{ik}|^{1+\delta}] < \infty$ for each $j, k \in \{1, \dots, d\}$ and some $\delta > 0$. Also, $\boldsymbol{\Sigma}_{xx} \equiv \text{plim} \mathbf{X}^\top \mathbf{X} / n$ exists, is finite and positive-definite with rank d . Finally, $\text{rank}[\mathbf{X}] = d$.

- AA3 allows regressors to be random (almost always the case with observational data).
- If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are iid then $\boldsymbol{\Sigma}_{xx} = E[\mathbf{x}_1 \mathbf{x}_1^\top]$.
- If $\mathbf{x}_1, \dots, \mathbf{x}_n$ are iid then AA3 guarantees the existence of $\boldsymbol{\Sigma}_{xx} \equiv \text{plim} \frac{1}{n} \mathbf{X}^\top \mathbf{X}$ via the Markov Law of Large Numbers.

AA4 $E[u_i | \mathbf{x}_i] = 0$ w.p. 1 for each i .

- AA2, AA4 imply that $E[\mathbf{y} | \mathbf{X}] = \mathbf{X} \boldsymbol{\beta}$ so $\boldsymbol{\beta}$ is intrinsically interesting.
- AA4 implies that $\text{Cov}[\mathbf{x}_i, u_i] = \mathbf{0}$ for each i , which suffices for OLS to be consistent.
- AA4 is a weak exogeneity assumption, i.e. u_i represents the total effect on y_i of all determinants of y_i that are uncorrelated with \mathbf{x}_i and on average u_i has no effect on y_i .
- Failure of AA4 is usually referred to as the presence of endogenous regressors. In this case, may not be consistent.
- AA4 does not imply that u_i, \mathbf{x}_i are independent.
- Strong exogeneity is the assumption that $E[u_i | \mathbf{X}] = 0$, i.e. $E[u_i | \mathbf{x}_j] = 0 \ \forall i, j \in \{1, \dots, n\}$. Under strong exogeneity, OLS is unbiased for $\boldsymbol{\beta}$ (refer to Hayashi Chapter 1).

AA5 $E[u_i^2 | \mathbf{x}_i] = \sigma_i^2$ for each $i \in \{1, \dots, n\}$.

$$\text{Also, } \boldsymbol{\Omega} \equiv E[\mathbf{u} \mathbf{u}^\top | \mathbf{X}] = \begin{bmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_n^2 \end{bmatrix} \text{ is positive definite, } \text{ i.e. } \sigma_i > 0 \ \forall i$$

$$\text{Finally, } E[|u_i^2|^{1+\delta}] < \infty \text{ for some } \delta > 0 \text{ and for all } i$$

- AA5 clearly allows the errors to exhibit conditional heteroskedasticity.

AA6 The matrix $\mathbf{S} \equiv \text{plim} \frac{1}{n} \sum_{i=1}^n u_i^2 \mathbf{x}_i \mathbf{x}_i^\top$ exists, is finite, positive definite and has rank d .

$$\text{In addition, } \exists \delta > 0 \text{ s.t. } E[|u_i^2 \mathbf{x}_{ij} \mathbf{x}_{ik}|^{1+\delta}] < \infty \text{ for each } i \text{ and each } j, k \in \{1, \dots, d\}$$

- AA6 guarantees that the asymptotic distribution of $\frac{1}{\sqrt{n}} \mathbf{X}^\top \mathbf{u}$ has a finite variance.

Technical aside:

a) Law of Large Numbers: We've seen the following.

Theorem (Kolmogorov LLN) Let $\{\mathbf{x}_i\}$ be iid. Let $\bar{\mathbf{x}}_n \equiv \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. If $E[\|\mathbf{x}_1\|] < \infty$, then

$$\bar{\mathbf{x}}_n \xrightarrow{\text{a.s.}} E[\mathbf{x}_n] = E[\mathbf{x}_1]$$

Theorem (Ergodic theorem) Let $\{\mathbf{x}_i\}$ be strictly stationary with an α -mixing base such that $E[\|\mathbf{x}_1\|] < \infty$, then

$$\bar{\mathbf{x}}_n \xrightarrow{\text{a.s.}} E[\mathbf{x}_n] = E[\mathbf{x}_1]$$

LLN for independent, but perhaps non-identically distributed processes

Theorem (Markov LLN) Let $\{x_i\}$ be independent random variables such that $\text{Var}[x_i] < \infty$ for each i . Let $\mu_i \equiv E[x_i]$, $\sigma_i^2 \equiv \text{Var}[x_i]$.

$$\text{If } \sum_{i=1}^{\infty} \frac{E[|x_i - \mu_i|^{1+\delta}]}{i^{1+\delta}} < \infty \text{ for some } \delta > 0, \quad \text{then } |\bar{x}_n - E[\bar{x}_n]| \xrightarrow{\text{a.s.}} 0$$

b) Central Limit Theorems: We've seen the following.

Theorem (Lindberg–Lévy CLT) Let $\{x_i\}$ be iid random variables with $\text{Var}[x_i] < \infty$, then

$$\frac{\bar{x}_n - E[\bar{x}_n]}{\sqrt{\text{Var}[\bar{x}_n]}} \xrightarrow{d} N(0,1)$$

Theorem (CLT for Martingale Difference Process) Let $\{x_i\}$ be strictly stationary, scalar valued MDP with α -mixing base. Suppose $E[x_i^2] < \infty$. Let $\sigma^2 \equiv E[x_1^2]$, then

$$\frac{\bar{x}_n}{\sigma/\sqrt{n}} \xrightarrow{d} N(0,1)$$

Now present a CLT for independent but not necessarily identically distributed processes.

Theorem (Lyapunov CLT) Let $\{x_i\}$ be independent random variables each with $\text{Var}[x_i] < \infty$. Let $\mu_i \equiv E[x_i]$, $\sigma_i^2 \equiv \text{Var}[x_i]$. Then

$$\text{if } \exists \delta > 0 \text{ such that } \lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n E[|x_i - \mu_i|^{2+\delta}]}{(\sum_{i=1}^n \sigma_i^2)^{\frac{2+\delta}{2}}} = 0, \quad \text{then } \frac{\bar{x}_n - E[\bar{x}_n]}{\sqrt{\text{Var}[\bar{x}_n]}} \xrightarrow{d} N(0,1)$$

c) Weighted Least Squares (WLS)

Rule of thumb: The existence of robust standard errors (e.g. heteroscedasticity-robust standard errors for OLS) usually implies that the existence of a more efficient estimator that avoids the need for robust standard errors. In the case of linear regression under AA1–AA6 generalized least squares (GLS) is more efficient than just OLS.

c) Weighted Least Squares (WLS) methodology

Suppose that the error variance $E[\mathbf{u}\mathbf{u}^T|\mathbf{X}] = \mathbf{\Omega}$, where $\mathbf{\Omega} \neq \sigma^2 \mathbf{I}_n$ and $\mathbf{\Omega}$ is nonsingular.

$$\text{Then } \mathbf{\Omega}^{-1/2} \text{ exists, and } \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \Leftrightarrow \mathbf{\Omega}^{-1/2}\mathbf{y} = \mathbf{\Omega}^{-1/2}\mathbf{X}\boldsymbol{\beta} + \mathbf{\Omega}^{-1/2}\mathbf{u}, \quad \dots (*)$$

Where

$$\text{Var}[\mathbf{\Omega}^{-1/2}\mathbf{u}|\mathbf{X}] = E[\mathbf{\Omega}^{-1/2}\mathbf{u}\mathbf{u}^T\mathbf{\Omega}^{-1/2}|\mathbf{X}] = \mathbf{I}_n$$

i.e. the errors in (*) are homoskedastic, and satisfy $E[\mathbf{\Omega}^{-1/2}\mathbf{u}|\mathbf{X}] = \mathbf{0}$ if $E[\mathbf{u}|\mathbf{X}] = \mathbf{0}$ is assumed. So by Gauss–Markov theorem, OLS applied to (*) is BLUE. The OLS estimator of $\boldsymbol{\beta}$ in (*) is the generalized least squares (GLS) estimator of $\boldsymbol{\beta}$. In particular,

$$\hat{\boldsymbol{\beta}}_{\text{GLS}} = (\mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{y}$$

Naturally, $\mathbf{\Omega}$ is unknown, so $\hat{\boldsymbol{\beta}}_{\text{GLS}}$ is infeasible. Instead, specify a parametric model for $\mathbf{\Omega}$, i.e. $\mathbf{\Omega} = \mathbf{\Omega}(\boldsymbol{\gamma})$, where $\mathbf{\Omega}(\cdot)$ is a known function and $\boldsymbol{\gamma}$ is a finite-dimensional vector. If the model $\mathbf{\Omega}(\boldsymbol{\gamma})$ is correctly specified and $\hat{\boldsymbol{\gamma}}$ is consistent for $\boldsymbol{\gamma}$, then $\hat{\mathbf{\Omega}} = \mathbf{\Omega}(\hat{\boldsymbol{\gamma}})$ is consistent for $\mathbf{\Omega}$.

e.g. suppose multiplicative heteroskedasticity, i.e.

$$E[u_i^2 | \mathbf{z}_i] = \exp(\mathbf{z}_i^T \boldsymbol{\gamma})$$

Where \mathbf{z}_i is a subvector of \mathbf{x}_i . Then $\boldsymbol{\gamma}$ is consistently estimated by OLS applied to a regression of $\log \hat{u}_i^2$ (where \hat{u}_i^2 is the i-th OLS residual) on \mathbf{z}_i ; or by nonlinear least squares applied to a regression of \hat{u}_i^2 on $\exp(\mathbf{z}_i^T \boldsymbol{\gamma})$.

Given an estimate $\hat{\mathbf{\Omega}} = \mathbf{\Omega}(\hat{\boldsymbol{\gamma}})$ of $\mathbf{\Omega}$ the feasible GLS (FGLS) estimator of $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}}_{\text{FGLS}} = (\mathbf{X}^T \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{\Omega}}^{-1} \mathbf{y}$$

Under assumptions AA1–AA6, if the model $\mathbf{\Omega}(\boldsymbol{\gamma})$ is correctly specified and $\hat{\boldsymbol{\gamma}}$ is consistent, then

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{FGLS}} - \boldsymbol{\beta}) \xrightarrow{d} N_d \left(0, \left(\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}^T \mathbf{\Omega}^{-1} \mathbf{X} \right)^{-1} \right), \quad \dots (**)$$

i.e. the same asymptotic variance as that of $\hat{\boldsymbol{\beta}}_{\text{GLS}}$, i.e. for large n , $\hat{\boldsymbol{\beta}}_{\text{FGLS}}$ is as efficient as $\hat{\boldsymbol{\beta}}_{\text{GLS}}$. Note that the FGLS estimated asymptotic variance matrix is

$$(\mathbf{X}^T \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1}$$

Compare to

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tilde{\mathbf{\Omega}}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}, \quad \text{where } \tilde{\mathbf{\Omega}} = \text{diag}[\hat{u}_i^2]$$

i.e. the estimated heteroscedasticity-robust asymptotic variance matrix for OLS (the FGLS estimated asymptotic variance implies standard errors that are no greater the corresponding OLS standard errors.). Naturally, (**) depends on $\mathbf{\Omega}(\boldsymbol{\gamma})$ being correctly specified. If $\mathbf{\Omega}(\boldsymbol{\gamma})$ is incorrectly specified, then $\hat{\boldsymbol{\beta}}_{\text{FGLS}}$ is still consistent, although $(\mathbf{X}^T \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1}$ will be inconsistent for the asymptotic variance.

A solution to possible misspecification of $\Omega(\gamma)$ involves a working variance matrix. Let $\Sigma \equiv \Sigma(\gamma)$ be a working variance matrix that may not equal the true variance matrix $\Omega \equiv E[\mathbf{u}\mathbf{u}^T|\mathbf{X}]$. Let $\hat{\Sigma} \equiv \Sigma(\hat{\gamma})$ where $\hat{\gamma}$ estimates γ . Then use WLS to β (with weighting matrix $\hat{\Sigma}^{-1}$), i.e.

$$\hat{\beta}_{\text{WLS}} = (\mathbf{X}^T \hat{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma}^{-1} \mathbf{y}$$

The estimated asymptotic variance matrix in this case is the heteroskedasticity-robust estimated asymptotic variance for OLS applied to

$$\hat{\Sigma}^{-1/2} \mathbf{y} = \hat{\Sigma}^{-1/2} \mathbf{X} \beta + \hat{\Sigma}^{-1/2} \mathbf{u}$$

i.e.

$$\begin{aligned} \text{Var}[\hat{\beta}_{\text{WLS}}] &= (\mathbf{X}^T \hat{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma}^{-1} \tilde{\Omega} \hat{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \hat{\Sigma}^{-1} \mathbf{X})^{-1} \\ \text{where } \tilde{\Omega} &= \text{diag}[(\hat{u}_i^*)^2] \\ \hat{u}_i^* &= y_i - \mathbf{x}_i^T \hat{\beta}_{\text{WLS}} \end{aligned}$$

In practice, working variance matrices typically involve only simple models for heteroskedasticity.

e.g. Numerical Example

$$\begin{aligned} \text{Suppose } y &= 1 + x + u \\ \text{where } u &= xv \\ x &\sim N(0, 5^2) \\ v &\sim N(0, 2^2), \quad \text{is independent of } x \end{aligned}$$

Then $\text{Var}[u|x] = x^2 \text{Var}[v] = 4x^2$, i.e. heteroskedastic. So OLS is inefficient, GLS is relatively efficient, while the efficiency of a working variance matrix approach is typically somewhere in between.

Consider $n=100$ simulated (iid) observations from this model.

- OLS generates a slope estimate of $\hat{\beta}_2 = 1.0495$. The heteroskedasticity-robust standard error is .3207 (correct). Compare with the homoscedasticity-only standard error .2016 (incorrect).

In fact, in this model, the heteroskedasticity-robust standard error for the slope is $\sqrt{12/n}$ when n is large; while the homoscedasticity only standard error for large n is $\sqrt{4/n}$.

- GLS for this model involves OLS applied to the transformed regression in which both sides of the original model are divided by $|x|$. Then the errors in the transformed model are homoskedastic; so the heteroskedasticity-robust and homoscedasticity-only standard errors should be close.

In fact, we see .208 (heteroskedasticity-robust) and .209 (homoscedasticity-only).

- Examples of the working variance matrix approach. Divided both sides of the original model by $\sqrt{|x|}$ (not the correct weights). The corresponding WLS estimator has heteroskedasticity-robust standard errors of .232.
- Summary: Compare standard errors for slope estimates in this model of .3207 (OLS), .232 (WLS based on a working variance matrix), .209 (GLS).
- Most applied work in practice does not involve weighted least squares techniques, perhaps for reasons of convenience.

d) Misspecification of AA1–AA6

Suppose at least one of AA1–AA6 is false.

(1) Then OLS might be inconsistent.

Recall: If $\mathbf{y}=\mathbf{X}\boldsymbol{\beta}+\mathbf{u}$ and if $\frac{1}{n}\mathbf{X}^\top\mathbf{u} \xrightarrow{p} \mathbf{0}$ then $\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}$ is consistent.

In fact, $\widehat{\boldsymbol{\beta}} = \boldsymbol{\beta} + \left(\frac{1}{n}\mathbf{X}^\top\mathbf{X}\right)^{-1}\frac{1}{n}\mathbf{X}^\top\mathbf{u}$ is consistent only if both $\mathbf{y}=\mathbf{X}\boldsymbol{\beta}+\mathbf{u}$ and $\text{plim}\frac{1}{n}\mathbf{X}^\top\mathbf{u} = \mathbf{0}$ are true. So $\widehat{\boldsymbol{\beta}}$ is inconsistent if $\mathbf{y}=\mathbf{X}\boldsymbol{\beta}+\mathbf{u}$ is false or if the regressors are asymptotically correlated with the error.

(2) If $\mathbf{y}=\mathbf{X}\boldsymbol{\beta}+\mathbf{u}$ is false, then OLS might not be useful as a predictor of \mathbf{y} , because in this case

$$E[\mathbf{y}|\mathbf{X}] = \mathbf{X}\boldsymbol{\beta}$$

is false.

e.g. Suppose $y_i=g(\mathbf{x}_i)+v_i$ for some nonlinear function $g(\cdot)$ and where $E[v_i|\mathbf{x}_i]=0$. Does OLS have a meaningful interpretation in this case? White (1980) showed that $E[y_i|\mathbf{x}_i]=g(\mathbf{x}_i)$ for some nonlinear function $g(\cdot)$, then OLS estimates $\boldsymbol{\beta} \in \mathbb{R}^d$ such that

$$\boldsymbol{\beta} = \underset{\mathbf{b} \in \mathbb{R}^d}{\text{argmin}} E[(g(\mathbf{x}) - \mathbf{x}^\top\mathbf{b})^2]$$

i.e. OLS generates a best linear predictor of \mathbf{y} with respect to mean-squared prediction error.

(Note: OLS cannot have a causal interpretation in this case.)

(3) If $\text{plim}\frac{1}{n}\mathbf{X}^\top\mathbf{u} = \mathbf{c} \neq \mathbf{0}$, i.e. the error and regressors are asymptotically correlated, then say that at least one element of \mathbf{x} is endogenous.

In this case, $E[u_i|\mathbf{x}_i]$ cannot be 0 and OLS is inconsistent.

Corresponding SAS Code

```
resetline;
ods html close;
ods graphics off;
ods listing;

data _01;
  do M=1 to 5;
    do N=1 to 100;
      X=5*rannor(1);
      V=2*rannor(2);
      U=X*V;
      Y=1+X+U;
      TX1=1/abs(X);
      TX2=X/abs(X);
      TY=Y/abs(X);
      UX1=1/sqrt(abs(X));
      UX2=X/sqrt(abs(X));
      UY=Y/sqrt(abs(X));
      output;
    end;
  end;
run;

proc reg;
  model Y=X/white;
  model TY=TX1 TX2/noint white;
  model UY=UX1 UX2/noint white;
  by M;
run;

quit;
```

(3) Suppose at least one of the chose regressors is endogenous.

Endogeneity of a regressor is usually defined to exist iff that regressor is not correlated with the error term.

This will be the case if AA4 fails, i.e. if $E[u_i|x_i]=0$ w.p. 1 is not true.

If a regressor is correlated with the error term (for whatever reason), then OLS is inconsistent.

(4) A classic cause of regressor endogeneity is due to one or more omitted variables.

e.g. Suppose $y_i = x_i^T \beta + z_i \alpha + v_i$ where $[x_i^T y_i z_i]^T$ are observed for all i and v_i is unobserved with $E[v_i|x_i, z_i]=0$ w.p. 1. Then OLS estimator of $[\beta^T \alpha]^T$ is consistent.

Now suppose y_i is regressed on x_i alone (so, z_i is omitted). In this case we are really dealing with the model

$$y_i = x_i^T \beta + u_i, \quad \text{where } u_i = z_i \alpha + v_i \text{ is unobserved}$$

If $\text{Cov}[x_i, z_i] \neq 0$, then the OLS estimator of β is inconsistent, with the inconsistency in this case called omitted variable bias (although it would be more accurate to call it omitted variable inconsistency).

e.g. Suppose

y_i = individual log average hourly earnings

x_i = vector of individual attributes (including years of schooling)

z_i = unobservable individual ability (i.e. IQ)

for a sample of n workers.

Generally assume z_i and years of schooling to be positively correlated; z_i and y_i to be positively correlated. If true, this assumption would make the OLS estimate of the coefficient vector attached to x_i in a regression of y_i on x_i inconsistent.

One can sometimes pin down the direction of inconsistency due to an omitted variable.

Let $z \equiv [z_1 \cdots z_n]^T$ denote the omitted observations on z_i . Let $v \equiv [v_1 \cdots v_n]^T$. Then the model for y is given by

$$y = X\beta + z\alpha + v$$

and the OLS “estimate” of β in a regression of y on X alone is

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y \\ &= (X^T X)^{-1} X^T (X\beta + z\alpha + v) \\ &= \beta + \left(\frac{1}{n} X^T X \right)^{-1} \frac{1}{n} X^T z \alpha + \left(\frac{1}{n} X^T X \right)^{-1} \frac{1}{n} X^T v, \quad \text{where } \underbrace{\frac{1}{n} X^T v \xrightarrow{p} 0}_{\text{since } \text{Cov}[x_i, v_i] = 0 \text{ is assumed}} \end{aligned}$$

If $\left(\frac{1}{n} X^T X \right)^{-1} \frac{1}{n} X^T z \xrightarrow{p} \delta \neq 0$, then $\hat{\beta} \xrightarrow{p} \beta + \delta \alpha$, which shows that the direction of inconsistency in $\hat{\beta}$ depends on the sign of each component of $\delta \alpha$.

e.g. Suppose for a sample of n workers

$y_i \equiv$ individual log average hourly earnings

$x_i \equiv$ individual years of schooling

$z_i \equiv$ unobserved individual ability

Then, if $E[v_i|x_i, z_i]=0$, then

$\beta =$ (causal) return to an additional year of schooling

$\alpha =$ (causal) return to ability

In this example, expect $\delta_j \alpha$ to be positive, so OLS applied to a regression of y_i on x_i will overestimate β in the limit.

(5) Suppose we have parameter heterogeneity. Suppose that instead of

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + u_i$$

We have instead

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}_i + u_i, \quad \dots (*)$$

when $E[u_i|\mathbf{x}_i]=0$.

In this case, the minimum MSPE (mean squared prediction error) predictor of y_i is

$$E[y_i|\mathbf{x}_i] = \mathbf{x}_i^T \boldsymbol{\beta}_i$$

which varies across i .

A random coefficients model is one where $\boldsymbol{\beta}_i$ are iid random vectors independent of \mathbf{x}_i .

Let $\boldsymbol{\beta} = E[\boldsymbol{\beta}_i]$. Then rewrite (*) as

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + [u_i + \mathbf{x}_i^T (\boldsymbol{\beta}_i - \boldsymbol{\beta})], \quad \dots (**)$$

If conditions AA1, AA3–AA6 hold, then OLS will be consistent for $\boldsymbol{\beta}$.

Special case of (**): A random effects model is one where $\mathbf{x}_i = [1 \ \tilde{\mathbf{x}}_i^T]^T$ and only the first component of $\boldsymbol{\beta}_i$ (i.e. the intercept) is random iid and independent of $\tilde{\mathbf{x}}_i$.

Also, OLS applied to (**) is inconsistent for $\boldsymbol{\beta}$ if $\boldsymbol{\beta}_i$ is correlated with \mathbf{x}_i .

A leading special case of this is the fixed effects model, i.e. $\mathbf{x}_i = [1 \ \tilde{\mathbf{x}}_i^T]^T$ and only the first component (i.e. the intercept) of $\boldsymbol{\beta}_i$ is random, but correlated with $\tilde{\mathbf{x}}_i$ (need panel data to handle this).

e) Instrumental Variables

Suppose $y_i = \mathbf{x}_i^T \boldsymbol{\beta} + u_i$ where (unfortunately)

$$\text{Cov}[\mathbf{x}_i, u_i] \neq \mathbf{0}$$

i.e. \mathbf{x}_i contains an endogenous component.

In this case, OLS is inconsistent. Can estimate $\boldsymbol{\beta}$ consistently if there exists a random m -vector of instruments \mathbf{z}_i where $m \geq d$ where

- | | |
|---|---|
| (1) $E[\mathbf{z}_i u_i] = \mathbf{0}$ | (otherwise \mathbf{z}_i is <u>invalid</u>) |
| (2) $E[\mathbf{z}_i \mathbf{x}_i^T] \neq \mathbf{0}$ | (otherwise \mathbf{z}_i is irrelevant) |
| (3) $\ E[\mathbf{z}_i \mathbf{x}_i^T]\ $ is large (rather than small) | (otherwise \mathbf{z}_i is weak) |

Properties (1), (2) (as discussed in Hayashi Ch. 3) are necessary for consistency. Property (3) has to do with controlling small-sample bias (more later).

The condition $m \geq d$ is the order condition for identification of $\boldsymbol{\beta}$.

Say that $\boldsymbol{\beta}$ is identified by \mathbf{z}_i if $m \geq d$

<u>overidentified</u>	$m > d$ (GMM, 2SLS)
<u>just identified</u>	$m = d$ (IV regression)
<u>not identified</u>	$m < d$

$$\text{Let } \mathbf{Z} \equiv [\mathbf{z}_1 \cdots \mathbf{z}_n]^T, \quad \text{suppose } m = d$$

Then, the instrumental variable (IV) estimator of $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}}_{IV} = (\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T \mathbf{y}$$

This is consistent because

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{IV} &= (\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T \mathbf{y} \\ &= (\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T (\mathbf{X} \boldsymbol{\beta} + \mathbf{u}) \\ &= \boldsymbol{\beta} + \left(\frac{1}{n} \mathbf{Z}^T \mathbf{X} \right)^{-1} \frac{1}{n} \mathbf{Z}^T \mathbf{u} \end{aligned}$$

Where

$$\begin{aligned} \frac{\mathbf{Z}^T \mathbf{X}}{n} &\xrightarrow{p} E[\mathbf{z}_1 \mathbf{x}_1^T] = \boldsymbol{\Sigma}_{zx} \neq \mathbf{0}, \quad \text{if } \mathbf{z}_1 \text{ is relevant} \\ \text{and } \frac{\mathbf{Z}^T \mathbf{u}}{n} &\xrightarrow{p} E[\mathbf{z}_1 u_1] = \mathbf{0}, \quad \text{if } \mathbf{z}_1 \text{ is valid} \end{aligned}$$

Can show that (under AA5, where $\boldsymbol{\Omega} \equiv E[\mathbf{u} \mathbf{u}^T | \mathbf{X}] = \text{diag}[\sigma_i^2 : i=1, \dots, n]$ and AA6) $\hat{\boldsymbol{\beta}}_{IV}$ is \sqrt{n} -consistent and asymptotically Normal with asymptotic variance consistency estimated as

$$\hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}_{IV}] \equiv (\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T \hat{\boldsymbol{\Omega}} \mathbf{Z} (\mathbf{Z}^T \mathbf{X})^{-1}$$

Where $\hat{\boldsymbol{\Omega}} = \text{diag}[\hat{u}_i^2 : i = 1, \dots, n]$, $\hat{u}_i = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_{IV}$.

In practice: $\hat{\boldsymbol{\beta}}_{IV}$ is often associated with large standard errors.

Now suppose $m \geq d$ (this allows for $m > d$). Consider the two-stage least squares (2SLS) estimator

$$\begin{aligned}\hat{\beta}_{2SLS} &= [\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{P}_Z \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_Z \mathbf{y}\end{aligned}$$

Where $\mathbf{P}_Z = \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$ is the projection matrix onto the column space of \mathbf{Z} . \mathbf{P}_Z is symmetric and idempotent, so

$$\hat{\beta}_{2SLS} = (\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top \mathbf{y}, \quad \text{where } \hat{\mathbf{X}} = \mathbf{P}_Z \mathbf{X}$$

Can show that $\hat{\beta}_{2SLS}$ is \sqrt{n} -consistent and asymptotically Normal with asymptotic variance consistently estimated by

$$\begin{aligned}\hat{\mathbf{V}}[\hat{\beta}_{2SLS}] &= n(\mathbf{X}^\top \mathbf{P}_Z \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \hat{\mathbf{S}}_n (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{X}^\top \mathbf{P}_Z \mathbf{X})^{-1} \\ \text{where } \hat{\mathbf{S}}_n &= \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \mathbf{z}_i \mathbf{z}_i^\top \\ \hat{u}_i &= y_i - \mathbf{x}_i^\top \hat{\beta}_{2SLS}, \quad (\text{note: need conditions AA5 and AA6})\end{aligned}$$

e.g. Suppose

$$\begin{aligned}y &= .5x + u, \quad \text{where } x = z + v \\ z &\sim N(2, 1^2) \\ \begin{bmatrix} u \\ v \end{bmatrix} &\sim N_2\left(\mathbf{0}, \begin{bmatrix} 1^2 & .8 \\ .8 & 1^2 \end{bmatrix}\right) \\ z &\text{ is independent of } [u \quad v]^\top\end{aligned}$$

We have a problem because

$$\text{Cov}[x, u] = \text{Cov}[z, u] + \text{Cov}[v, u] = \text{Cov}[v, u] = .8$$

So, OLS applied to a regression of y on x is inconsistent for $.5$.

On the other hand, IV estimation with instrument z is consistent because z is valid ($E[zu] = \text{Cov}[z, u] + E[z]E[u] = 0 + 0 = 0$) and relevant ($E[zx] = \text{Cov}[z, x] + E[z]E[x] = 1 + 4 = 5 \neq 0$).

In addition, certain functions of z (e.g. z^3) can be used as instruments. Consider $n=10,000$ simulated samples from this model.

The OLS estimate (robust standard error) is $.902$ (.006).

The IV estimate using z is $.510$ (.010).

The 2SLS estimate using z is $.510$ (.014).

The IV estimate using z^3 is $.509$ (.012).

Obviously the IV and 2SLS standard errors are larger than the OLS standard error; in this case this is compensated by much smaller bias.

Corresponding SAS Code

```
resetline;
ods html close;
ods graphics off;
ods listing;

/*
proc iml;
    VAR={1 .8, .8 1};
    HALF=root(VAR);
    print HALF;
quit;
*/

data _01;
    do N=1 to 10000;
        Z=2+rannor(1);
        Z3=Z*3;
        U=rannor(2);
        V=.8*U+.6*rannor(3);
        X=Z+V;
        Y=.5*X+U;
        output;
    end;
run;

proc reg;
    model Y=X/white;
run;

proc iml;
    use _01;
    read all var{Y} into Y;
    read all var{X} into X;
    X=j(nrow(X),1,1)||X;
    read all var{Z} into Z;
    Z=j(nrow(Z),1,1)||Z;
    IV=inv(Z`*X)*Z`*Y;
    E1=Y-X*IV;
    VIV=inv(Z`*X)*Z`*diag(E1##2)*Z*inv(Z`*X);
    SEIV=sqrt(vecdiag(VIV));
    PZ=Z*inv(Z`*Z)*Z`;
    TSLS=inv(X`*PZ*X)*X`*PZ*Y;
    E2=Y-X*TSLS;
    S=0;
    do N=1 to 10000;
        S=S+E2[N]**2*Z[N,]*Z[N,];
    end;
    VTSL=inv(X`*PZ*X)*X`*Z*inv(Z`*Z)*S*inv(Z`*Z)*Z`*X*inv(X`*PZ*X);
    SETSLS=sqrt(vecdiag(VTSL));
    read all var{Z3} into Z3;
    Z3=j(nrow(Z3),1,1)||Z3;
    IV3=inv(Z3`*X)*Z3`*Y;
    E3=Y-X*IV3;
    VIV3=inv(Z3`*X)*Z3`*diag(E1##2)*Z3*inv(Z3`*X);
    SEIV3=sqrt(vecdiag(VIV3));
    print IV[format=8.4]
           SEIV[format=8.4],
           TSLS[format=8.4]
           SETSLS[format=8.4],
           IV3[format=8.4]
           SEIV3[format=8.4];
quit;

proc model;
    Y=B0+B1*X;
    instruments Z;
    fit Y/2sls hccme=0;
run;

quit;
```

Recall: IV/2SLS tends to have smaller bias than OLS, but at the cost of larger standard errors.

Today: If the chosen instruments are weak, then IV/2SLS can be as biased (or even more biased) than OLS while containing to be relatively inefficient.

Weak instruments are associated with

$$\|E[\mathbf{z}_i \mathbf{x}_i^T]\|$$

being “small.”

(Instrument weakness causes finite sample bias)

e.g. Consider

$$y = \beta_1 x_1 + \mathbf{x}_2^T \boldsymbol{\beta}_2 + u, \quad \dots (*)$$

where $\text{Cov}[x_1, u] \neq 0$
and $\text{Cov}[\mathbf{x}_2, u] = \mathbf{0}$

Assume that the instrument \mathbf{z} includes \mathbf{x}_2 and at least one other exogenous variable.

Intuition: \mathbf{z} is weak if the R^2 from a regression of x_1 on \mathbf{z} is “low,” i.e. \mathbf{z} are not strongly correlated.

But: Really only need x_1 to be strongly correlated with the elements of \mathbf{z} not in \mathbf{x}_2 .

Bound, Jaeger and Baker (1995, JASA): Proposal for a partial R^2 , i.e. the R^2 from the regression

$$x_1 - \hat{x}_1 = (\mathbf{z} - \hat{\mathbf{z}})^T \boldsymbol{\gamma} + v_1, \quad \dots (**)$$

Where \hat{x}_1 is the LS predicted value from a regression of x_1 on \mathbf{x}_2 and $\hat{\mathbf{z}}$ is the vector of LS predicted values from regressions of each element of \mathbf{z} on \mathbf{x}_2 .

Case with more than one endogenous regressor, i.e. the case where \mathbf{x}_2 in (*) contains an endogenous regressor.

In this case, can diagnose instrument weakness via the proposal of Shea (1997, RES): Use the R^2 from the following regression

$$x_1 - \hat{x}_1 = (\tilde{x}_1 - \tilde{\tilde{x}}_1)\eta + v_2$$

Where \hat{x}_1 is the LS predicted value from a regression of x_1 on \mathbf{x}_2 and \tilde{x}_1 is the LS predicted value from a regression of x_1 on \mathbf{z} ; $\tilde{\tilde{x}}_1$ is the LS predicted value from a regression of \tilde{x}_1 on $\tilde{\mathbf{x}}_2$; where $\tilde{\mathbf{x}}_2$ is the LS predicted value from a regression of \mathbf{x}_2 on \mathbf{z} .

Another method of diagnosing instrument weakness uses the F-statistic for the hypothesis that all coefficients are zero in a regression of an endogenous regressor on instruments.

e.g. Go back to (*). Suppose $\mathbf{z} = [\mathbf{z}_1^T \mathbf{z}_2^T]^T$ where \mathbf{z}_1 are exogenous variables not in \mathbf{x}_2 . Then can use the F-statistic for $H_0: \boldsymbol{\pi}_1 = \mathbf{0}$ in the “reduced-form” regression

$$x_1 = \mathbf{z}_1^T \boldsymbol{\pi}_1 + \mathbf{z}_2^T \boldsymbol{\pi}_2 + v_3$$

Rule of thumb: \mathbf{z} is “very weak” if this F-statistic is ≤ 5 ; \mathbf{z} is “potentially weak” if the F-statistic is (5,10).

Next: IV/2SLS is inconsistent if the chosen instrument is invalid, i.e. if \mathbf{z} is correlated with \mathbf{u} .

Deciding on questions of instrument validity is usually subjective, but note that instrument validity is closely tied to a model classification.

e.g.

$$\begin{aligned} y &= x_1\beta_1 + x_2\beta_2 + u \\ \text{where } \text{Cov}[x_1, u] &= 0 \\ \text{but } \text{Cov}[x_2, u] &\neq 0 \\ \text{and } E[x_1] &= E[x_2] = 0 \end{aligned}$$

Then OLS is inconsistent. “Plausible” that x_1^2 can serve as an instrument. This may not work, however, because it assumes that $E[y|x_1, x_2]$ is linear but not quadratic in x_1 .

Better not to tie instrument validity to a restriction on the function form of $E[y|x_1, x_2]$.

Next: The bias due to very slight instrument endogeneity is much worse if the instrument is also weak.

e.g.

$$y = x\beta + u$$

Suppose z is a candidate instrument. Assume $E[x]=E[z]=E[u]=0$ and $\text{Var}[x]=\text{Var}[z]=\text{Var}[u]=1^2$.

Suppose we observe (x_i, y_i, z_i) ($i=1, \dots, n$) iid. We have

$$\begin{aligned} \hat{\beta}_{IV} &= \frac{\sum_{i=1}^n z_i y_i}{\sum_{i=1}^n z_i x_i}, & \hat{\beta}_{OLS} &= \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \\ &\xrightarrow{p} \frac{E[z_1 y_1]}{E[z_1 x_1]} = \beta + \frac{E[z_1 u_1]}{E[z_1 x_1]} \end{aligned}$$

and

$$\hat{\beta}_{OLS} \xrightarrow{p} \frac{E[x_1 y_1]}{E[x_1^2]} = E[x_1 y_1] = \beta + E[x_1 u_1]$$

So

$$\frac{\text{plim } \hat{\beta}_{IV} - \beta}{\text{plim } \hat{\beta}_{OLS} - \beta} = \frac{E[z_1 u_1]}{E[z_1 x_1]} \frac{1}{E[x_1 u_1]} = \frac{\text{Corr}[z_1, u_1]}{\text{Corr}[x_1, u_1]} \frac{1}{\text{Corr}[z_1, x_1]}$$

i.e. if z is invalid and weak, then $\hat{\beta}_{IV}$ will be even more inconsistent for β than $\hat{\beta}_{OLS}$. Then

$$\begin{aligned} |\text{plim } \hat{\beta}_{IV} - \beta| &> |\text{plim } \hat{\beta}_{OLS} - \beta| \\ \text{if } |\text{Corr}[z_1, u_1]| &> .10 |\text{Corr}[x_1, u_1]| \end{aligned}$$

This shows that only a slight amount of instrument endogeneity suffices to make a weak instrument generate a large degree of inconsistency in the corresponding IV estimate.

Also: IV/2SLS become even more imprecise when the chosen instruments are weak.

e.g.

$$y = x\beta + u, \quad \text{Corr}[x, u] \neq 0$$

Suppose z is a valid instrument, and we observe (x_i, y_i, z_i) ($i=1, \dots, n$) iid. Assume that the corresponding errors u_1, \dots, u_n are iid mean zero with variance σ^2 .

In this case, can use a “homoscedasticity-only” estimator of the $\hat{\beta}_{IV}$ asymptotic variance

$$\begin{aligned} \hat{V}[\hat{\beta}_{IV}] &= s^2 (\mathbf{x}^T \mathbf{z})^{-1} \mathbf{z}^T \mathbf{z} (\mathbf{z}^T \mathbf{x})^{-1} \\ &= \frac{s^2}{\mathbf{x}^T \mathbf{x}} \frac{(\mathbf{z}^T \mathbf{z})(\mathbf{x}^T \mathbf{x})}{(\mathbf{z}^T \mathbf{x})^2} = \hat{V}[\hat{\beta}_{OLS}] \frac{1}{R_{x|z}^2} + o_p(1) \end{aligned}$$

where $R_{x|z}^2$ is the squared sample correlation between x and z and where $\hat{V}[\hat{\beta}_{OLS}]$ is the homoscedasticity-only OLS asymptotic variance estimator.

Therefore, the relative inefficiency of $\hat{\beta}_{IV}$ (relative to $\hat{\beta}_{OLS}$) is made worse when $R_{x|z}^2$ is small.

e.g. Suppose

$$y = x_1 \beta_1 + \mathbf{x}_2^T \boldsymbol{\beta}_2 + u$$

where \mathbf{x}_1 is endogenous and \mathbf{x}_2 is exogenous.

Suppose \mathbf{z} is a valid instrument that includes \mathbf{x}_2 and at least one additional exogenous variable.

Let R_p^2 be the R^2 from the regression

$$\mathbf{x}_1 - \hat{\mathbf{x}}_1 = (\mathbf{z} - \hat{\mathbf{z}})^T \boldsymbol{\gamma} + v$$

where $\hat{\mathbf{x}}_1$ is the LS fitted value of x_1 from a regression of x_1 on \mathbf{x}_2 , and $\hat{\mathbf{z}}$ is the vector of LS fitted values from regressions of each component of \mathbf{z} on \mathbf{x}_2 .

Let $\hat{\beta}_{1,2SLS}$ be the 2SLS estimator of β_1 using \mathbf{z} .

Let $\hat{\beta}_{1,OLS}$ be the OLS estimator of β_1 . Can show

$$SE[\hat{\beta}_{1,2SLS}] = \frac{SE[\hat{\beta}_{1,OLS}]}{R_p}$$

So, $SE[\hat{\beta}_{1,2SLS}]$ will generally be larger than $SE[\hat{\beta}_{1,OLS}]$.

Also, instrumental weakness affects IV/2SLS standard errors only for estimated coefficients attached to endogenous variables.

Instrument weakness inflates IV/2SLS standard errors via R_p^2 (i.e. partial R^2) and not through the R^2 on the first-stage regression.

(So, 1) compare SE_{2SLS} and SE_{OLS} and 2) if there is a significant gap between them, have to think about instrument weakness)

(This is because that correlation between an endogenous regressor and exogenous regressors is unimportant; want correlation between an endogenous regressor and instruments that are not also exogenous regressors)

2SLS standard errors on estimated coefficients attached to endogenous regressors that are much larger than the corresponding OLS standard errors indicate instrument weakness.

Relationship between instrument weakness and the finite-sample bias of IV or 2SLS estimators

Suppose \mathbf{Z} is an $(n \times m)$ -matrix of valid instruments.

Suppose $m=d=\#$ of regressors.

Simulations have shown that

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

can be highly biased in “small” samples (even for $n \geq 100,000$) despite being consistent.

In addition, the estimated asymptotic standard errors of $\hat{\beta}_{IV}$ can be highly biased even in “small” samples.

First consider that $E[\hat{\beta}_{IV}]$ is basically intractable.

$$\begin{aligned} E[\hat{\beta}_{IV}] &= E[(\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T \mathbf{y}] \\ &= \beta + E[(\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T \mathbf{u}] \\ &= \beta + E[(\mathbf{Z}^T \mathbf{X})^{-1} \mathbf{Z}^T E[\mathbf{u}|\mathbf{Z}, \mathbf{X}]] \end{aligned}$$

Therefore, $\hat{\beta}_{IV}$ is unbiased if $E[\mathbf{u}|\mathbf{Z}, \mathbf{X}] = \mathbf{0}$ w.p. 1.

But, if $E[\mathbf{u}|\mathbf{Z}, \mathbf{X}] = \mathbf{0}$ a.s. then $E[\mathbf{u}|\mathbf{X}] = \mathbf{0}$ a.s., which implies that \mathbf{X} is exogenous, i.e. $\hat{\beta}_{IV}$ is unnecessary. Therefore, the assumption that

$$E[\mathbf{u}|\mathbf{Z}, \mathbf{X}] = \mathbf{0}, \quad \text{w.p. 1}$$

is much too strong in this case. Therefore, we expect that

$$E[\hat{\beta}_{IV}] \neq \beta$$

in general.

Research into the value of

$$|E[\hat{\beta}_{IV}] - \beta|, \quad \text{or } |E[\hat{\beta}_{2SLS}] - \beta|$$

Take the case $d=1$, i.e. just one regressor, which is endogenous. Consider the reduced-form regression

$$x = \mathbf{z}^T \boldsymbol{\pi} + v, \quad \text{where } E[v] = 0 \text{ and } \text{Var}[v] = \sigma_v^2$$

Define the concentration parameter

$$\tau^2 = \frac{\boldsymbol{\pi}^T \mathbf{Z}^T \mathbf{Z} \boldsymbol{\pi}}{\sigma_v^2}$$

i.e. the signal-to-noise ratio in the reduced form regression.

Can show that $|E[\hat{\beta}_{2SLS}] - \beta|$ is decreasing in τ^2 , so having weak instruments tends to increase bias.

Note: (τ^2/m) is approximately the F-statistic (in fact the F-statistic minus one) for a test of $H_0: \boldsymbol{\pi} = \mathbf{0}$.

So: A test for finite-sample bias/instrument weakness can be based on (τ^2/m) .

In fact: For the case $m=d=1$, then $1/F$ is approximately unbiased for

$$\frac{|E[\hat{\beta}_{IV}] - \beta|}{|E[\hat{\beta}_{2SLS}] - \beta|}$$

(Straiger and Stock, 1997, EMA)

In general, the extent to which $\hat{\beta}_{IV}$ and $\hat{\beta}_{2SLS}$ is biased relative to $\hat{\beta}_{OLS}$ depends on

- The # of endogenous regressors
- M

Straiger and Stock's simulations proposed the rule of thumb.

- "F>10" to keep the bias of $\hat{\beta}_{IV}$ or $\hat{\beta}_{2SLS}$ to no more than .10 of $\|E[\hat{\beta}_{OLS}] - \beta\|$

Empirical example

Kling (2001, JBES) considered the causal effect of schooling on earnings. Kling used college proximity (to college) as an instrument for schooling.

The model is

$$\log W_i = \beta_0 + \beta_1 S_i + \beta_2 E_i + \beta_3 E_i^2 + \mathbf{x}_{2i}^T \boldsymbol{\gamma} + u_i, \quad i = 1, \dots, 3010$$

where $W_i \equiv$ average hourly earnings
 $S_i \equiv$ years of schooling
 $E_i \equiv \text{Age}_i - S_i - 6$
 $\mathbf{x}_{2i} \equiv$ vector of 26 additional variables

(Data were from the NLSY and consisted of observations on men aged 24–34 in 1976)
 Assume

- S_i, E_i, E_i^2 are endogenous because of unobservable ability
- Because of this, require at least 3 instruments not present in \mathbf{x}_{2i}
- Kling used three excluded instruments

$\text{Col4}_i \equiv$ indicator for proximity to a 4 year college
 $\text{Age}_i, \text{Age}_i^2$ in years and years squared

- The overall instrument vector was completed with a constant and \mathbf{x}_2 , i.e. $m=d=30$
- Kling tested for homoscedasticity and did not reject, so used homoscedasticity-only standard errors
- The OLS coefficient on S_i is .073 (.004)

Interpretation: An additional year of schooling causes wages to rise on average by $7.6\% = \exp(.073) - 1$

- The IV coefficient on S_i is .132 (.049)
- Shea's partial R^2 on schooling is $R_p^2 = .0064$, which implies that the IV standard error on schooling should be $1/.08 = 12.5$ times the OLS standard error, which is a close to what we observe
- Now consider the reduced-form regression of S_i on all instruments

Get an F-statistic of 8.07, which suggests that the finite-sample bias of the IV estimator is between 10 and 20% of that of OLS.

Summary: Working with instruments is inherently difficult. Ideally, one should avoid using instruments by exploiting “better” data from randomized controlled experiments or from a natural experiment

Nonlinear estimators

Suppose we observe

$$\{[\mathbf{x}_i^\top \ y_i]^\top : i = 1, \dots, n\}$$

A linear estimator is one that is linear in y_i .

e.g. OLS

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \sum_{i=1}^n \mathbf{x}_i y_i$$

i.e. weighted sum of the y_i s.

e.g. 2SLS

$$\hat{\boldsymbol{\beta}}_{2\text{SLS}} = (\mathbf{X}^\top \mathbf{P}_Z \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_Z \mathbf{y}$$

i.e. weighted sum of the y_i s.

A nonlinear estimator is one that is not linear in the y_i s. These are often associated with

- Estimates of conditional mean function parameters in conditional means that are not linear in parameters
- Estimators of parameters defined implicitly as solutions to optimization problems
- “Limited” dependent variable models, e.g. models in which the y_i s are generated via censoring or truncation

We will discuss

- 1) Extremum (M-) estimators

(i.e. estimators defined to solve an optimization problem)

- 2) Estimating-equations (Z-) estimators

(i.e. estimators defined to solve a system of estimating equations)

- 3) Maximum-likelihood estimators (MLEs) and quasi-maximum-likelihood estimators (QMLEs)
- 4) Nonlinear least squares estimators

Note: These four categories overlap and include linear estimators

e.g. OLS is an extremum estimator since

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \underset{\mathbf{b} \in \mathbb{R}^d}{\operatorname{argmin}} (\mathbf{y} - \mathbf{X}\mathbf{b})^\top (\mathbf{y} - \mathbf{X}\mathbf{b})$$

$\hat{\boldsymbol{\beta}}_{\text{OLS}}$ is also an estimating equations estimator because it is a solution to the system

$$\mathbf{X}^\top (\mathbf{X}\mathbf{b} - \mathbf{y}) = \mathbf{0}$$

$\hat{\boldsymbol{\beta}}_{\text{OLS}}$ is an MLE of $\boldsymbol{\beta}$ in the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \text{where } \mathbf{u}|\mathbf{X} \sim N(0, \sigma^2 \mathbf{I}_n)$$

e.g. Poisson regression

Observe $(\mathbf{x}_i^\top \mathbf{y}_i)^\top$ ($i=1, \dots, n$).

Where for each i the conditional distribution of y_i given \mathbf{x}_i is $\text{Poisson}(\exp(\mathbf{x}_i^\top \boldsymbol{\beta}))$, i.e. the conditional density of \mathbf{y}_i given \mathbf{x}_i is

$$f(y_i | \mathbf{x}_i; \boldsymbol{\beta}) = \begin{cases} \frac{\exp(-\exp(\mathbf{x}_i^\top \boldsymbol{\beta})) (\exp(\mathbf{x}_i^\top \boldsymbol{\beta}))^{y_i}}{y_i!}, & y_i = 0, 1, 2, \dots \\ 0, & \text{otherwise} \end{cases}$$

Assuming iid observations, the MLE of $\boldsymbol{\beta}$ maximizes

$$\prod_{i=1}^n f(y_i | \mathbf{x}_i; \boldsymbol{\beta}), \quad \text{as a function of } \boldsymbol{\beta}$$

i.e. the MLE of $\boldsymbol{\beta}$ maximizes

$$\begin{aligned} \widehat{Q}_n(\boldsymbol{\beta}) &\equiv \frac{1}{n} \log \prod_{i=1}^n f(y_i | \mathbf{x}_i; \boldsymbol{\beta}) \\ &= \frac{1}{n} \sum_{i=1}^n (-\exp(\mathbf{x}_i^\top \boldsymbol{\beta}) + y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log y_i!) \end{aligned}$$

i.e. the MLE of $\boldsymbol{\beta}$ solves the FOCs

$$\left. \frac{\partial}{\partial \boldsymbol{\beta}^\top} \widehat{Q}_n(\boldsymbol{\beta}) \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_n} = \mathbf{0} \Leftrightarrow \frac{1}{n} \sum_{i=1}^n (y_i - \exp(\mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}_n)) \mathbf{x}_i = \mathbf{0}$$

(No closed form expression for $\widehat{\boldsymbol{\beta}}_n$)

Extremum estimators: In general, an extremum estimator/M-estimator $\widehat{\boldsymbol{\theta}}_n$ of a parameter $\boldsymbol{\theta}_0 \in \mathbb{R}^d$ maximizes an objective function $\widehat{Q}_n(\boldsymbol{\theta}) \equiv \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, y_i, \boldsymbol{\theta})$ for a scalar-valued contrast function $g(\cdot, \cdot, \cdot)$.

Different extremum estimators correspond to different contrast functions (e.g. MLE $g(\cdot, \cdot, \cdot)$ is a log density) (e.g. OLS $g(\cdot, \cdot, \cdot)$ is the squared residual).

Consider an extremum estimator $\hat{\theta}_n$ of a parameter θ_0 , i.e.

$$\hat{\theta}_n = \operatorname{argmax}_{\theta} \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, y_i, \theta)$$

e.g. MLE of β in a Poisson regression

$$\begin{aligned} y|\mathbf{x} &\sim \text{Poisson}(e^{\mathbf{x}^\top \beta}) \\ \Rightarrow g(\mathbf{x}_i, y_i, \beta) &= -e^{\mathbf{x}_i^\top \beta} + y_i \mathbf{x}_i^\top \beta - \log y_i! \end{aligned}$$

e.g. OLS estimator of β in a linear model with

$$\begin{aligned} E[y|\mathbf{x}] &= \mathbf{x}^\top \beta \\ \Rightarrow g(\mathbf{x}_i, y_i, \beta) &= (y_i - \mathbf{x}_i^\top \beta)^2 \end{aligned}$$

We usually work with M-estimators $\hat{\theta}_n$ that solve systems of first-order conditions

$$\left. \frac{\partial}{\partial \theta} \hat{Q}_n(\theta) \right|_{\theta=\hat{\theta}_n} = \mathbf{0}$$

i.e. M-estimators that are also Z-estimators

An M-estimator is worthwhile if it is both consistent ($\hat{\theta}_n \xrightarrow{p} \theta_0$) and if it has a “tractable” asymptotic distribution, i.e. if there is an increasing sequence of constants $\{a_n\}$ s.t. $a_n(\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathbf{L}(\theta_0)$, where $\mathbf{L}(\theta_0)$ is a random vector whose distribution is a known function of θ_0 . This implies that if n is large, then $\mathbf{L}(\theta_0)$ has a distribution approximately equal to that of $\mathbf{L}(\hat{\theta}_n)$; i.e. tests/confidence regions regarding θ_0 can be based off the (known) distribution of $\mathbf{L}(\hat{\theta}_n)$.

Usual case in this course

$$a_n = \sqrt{n}, \quad \mathbf{L}(\theta_0) \sim N_d(\mathbf{0}, \mathbf{V}(\theta_0))$$

In particular, if $\hat{\theta}_n$ satisfies

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} g(\mathbf{x}_i, y_i, \theta) \Big|_{\theta=\hat{\theta}_n} = \mathbf{0}$$

then a Taylor expansion (assuming that it is feasible) shows that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = - \left(\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta^\top} g(\mathbf{x}_i, y_i, \theta) \Big|_{\theta=\bar{\theta}_n} \right)^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \theta} g(\mathbf{x}_i, y_i, \theta) \Big|_{\theta=\theta_0}$$

where $\bar{\theta}_n$ is a point on the line segment connecting $\hat{\theta}_n$ and θ_0 .

Therefore

$$\text{if } \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} g(\mathbf{x}_i, y_i, \boldsymbol{\theta}_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{B}_0)$$

where \mathbf{B}_0 is finite and non-singular and

$$\text{if } \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} g(\mathbf{x}_i, y_i, \boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n} \xrightarrow{p} \mathbf{A}_0$$

where \mathbf{A}_0 is finite and non-singular, then

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

by Slutsky's theorem.

Suppose $\hat{\mathbf{A}}_n \xrightarrow{p} \mathbf{A}_0$, $\hat{\mathbf{B}}_n \xrightarrow{p} \mathbf{B}_0$, i.e. $\mathbf{A}_0, \mathbf{B}_0$ can be consistently estimated. In this case, the asymptotic variance $\frac{1}{n} \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1}$ of $\hat{\boldsymbol{\theta}}_n$ can be consistently estimated as

$$\hat{\mathbf{V}}_n[\hat{\boldsymbol{\theta}}_n] \equiv \frac{1}{n} \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n \hat{\mathbf{A}}_n^{-1}$$

e.g. Suppose we observe $[\mathbf{x}_i^\top \ y_i]^\top$ ($i=1, \dots, n$) iid with $y_i | \mathbf{x}_i \sim \text{Poisson}(e^{\mathbf{x}_i^\top \boldsymbol{\beta}})$.

Then the Poisson MLE $\hat{\boldsymbol{\beta}}_n$ is given by

$$\hat{\boldsymbol{\beta}}_n = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, y_i, \boldsymbol{\beta})$$

where

$$g(\mathbf{x}_i, y_i, \boldsymbol{\beta}) = -e^{\mathbf{x}_i^\top \boldsymbol{\beta}} + y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log y_i!$$

So

$$\frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}_i, y_i, \boldsymbol{\beta}) = (-e^{\mathbf{x}_i^\top \boldsymbol{\beta}} + y_i) \mathbf{x}_i$$

and

$$\frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} g(\mathbf{x}_i, y_i, \boldsymbol{\beta}) = -e^{\mathbf{x}_i^\top \boldsymbol{\beta}} \mathbf{x}_i \mathbf{x}_i^\top$$

Therefore

$$\begin{aligned} \mathbf{A}_0 &= -\operatorname{plim} \frac{1}{n} \sum_{i=1}^n e^{\mathbf{x}_i^\top \boldsymbol{\beta}_0} \mathbf{x}_i \mathbf{x}_i^\top \\ \mathbf{B}_0 &= \operatorname{plim} \frac{1}{n} \sum_{i=1}^n (y_i - e^{\mathbf{x}_i^\top \boldsymbol{\beta}_0})^2 \mathbf{x}_i \mathbf{x}_i^\top \end{aligned}$$

The asymptotic variance of the Poisson MLE is consistently estimated as

$$\hat{\mathbf{V}}_n[\hat{\boldsymbol{\beta}}_n] = \frac{1}{n} \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n \hat{\mathbf{A}}_n^{-1}$$

where

$$\begin{aligned} \hat{\mathbf{A}}_n &= \frac{1}{n} \sum_{i=1}^n e^{\mathbf{x}_i^\top \boldsymbol{\beta}} \mathbf{x}_i \mathbf{x}_i^\top \\ \hat{\mathbf{B}}_n &= \frac{1}{n} \sum_{i=1}^n \left(y_i - e^{\mathbf{x}_i^\top \boldsymbol{\beta}} \right)^2 \mathbf{x}_i \mathbf{x}_i^\top \end{aligned}$$

Note: If $y_i | \mathbf{x}_i \sim \text{Poisson}(e^{\mathbf{x}_i^\top \boldsymbol{\beta}_0})$ then $E[y_i | \mathbf{x}_i] = \text{Var}[y_i | \mathbf{x}_i] = e^{\mathbf{x}_i^\top \boldsymbol{\beta}_0}$

So

$$\begin{aligned} E \left[\frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}, y, \boldsymbol{\beta}_0) \middle| \mathbf{x} \right] &= \mathbf{0} \\ \text{Var} \left[\frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}, y, \boldsymbol{\beta}_0) \middle| \mathbf{x} \right] &= \text{Var}[y | \mathbf{x}] \mathbf{x} \mathbf{x}^\top = e^{\mathbf{x}^\top \boldsymbol{\beta}} \mathbf{x} \mathbf{x}^\top \end{aligned}$$

which implies that

$$\mathbf{B}_0 = E[e^{\mathbf{x}^\top \boldsymbol{\beta}_0} \mathbf{x} \mathbf{x}^\top] = -\mathbf{A}_0$$

In this case, the asymptotic variance of the Poisson MLE can also be estimated as

$$\hat{\mathbf{V}}_n[\hat{\boldsymbol{\beta}}_n] = -\frac{1}{n} \hat{\mathbf{A}}_n^{-1}$$

which implies smaller standard errors than those available from $\hat{\mathbf{V}}_n[\hat{\boldsymbol{\beta}}_n]$. More popular in this case to use $\hat{\mathbf{V}}_n[\hat{\boldsymbol{\beta}}_n]$ because of its robustness against deviations from the assumption $y | \mathbf{x} \sim \text{Poisson}(e^{\mathbf{x}^\top \boldsymbol{\beta}_0})$.

Interpretation of coefficients in nonlinear regression models

Suppose $E[y | \mathbf{x}] = g(\mathbf{x}, \boldsymbol{\beta}_0)$ w.p. 1. Basically care about marginal effects, i.e. the vector

$$\frac{\partial}{\partial \mathbf{x}} E[y | \mathbf{x}]$$

as a function of \mathbf{x} . If $E[y | \mathbf{x}] = \mathbf{x}^\top \boldsymbol{\beta}$ (i.e. linear in $\boldsymbol{\beta}$) then the vector of marginal effect is constant in \mathbf{x} and in fact is $\boldsymbol{\beta}$.

In nonlinear models, $\frac{\partial}{\partial \mathbf{x}} E[y | \mathbf{x}]$ varies with \mathbf{x} . In this case, not usually easy to interpret $\boldsymbol{\beta}_0$.
e.g. $E[y | \mathbf{x}] = e^{\mathbf{x}^\top \boldsymbol{\beta}_0}$

$$\frac{\partial}{\partial \mathbf{x}} E[y | \mathbf{x}] = e^{\mathbf{x}^\top \boldsymbol{\beta}_0} \boldsymbol{\beta}_0$$

which varies clearly with \mathbf{x} .

In nonlinear models, the distribution of $\frac{\partial}{\partial \mathbf{x}} E[y|\mathbf{x}]$ gives a comprehensive picture of the relation between \mathbf{x} and y .

But: Difficult to interpret when \mathbf{x} is relatively high dimensional.

Instead, popular to focus on

(1)

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \mathbf{x}_i} E[y_i | \mathbf{x}_i]$$

i.e. vector of average marginal responses

(2)

$$\left. \frac{\partial}{\partial \mathbf{x}} E[y|\mathbf{x}] \right|_{\mathbf{x}=\bar{\mathbf{x}}}$$

i.e. marginal response for the “average individual”

(3)

$$\left. \frac{\partial}{\partial \mathbf{x}} E[y|\mathbf{x}] \right|_{\mathbf{x}=\mathbf{x}^*}$$

i.e. marginal response for an individual with characteristics \mathbf{x}^*

Note:

- (1)–(3) are all β_0 when $E[y|\mathbf{x}] = \mathbf{x}^T \beta_0$
- (1) is probably most relevant for policy analysis, followed closely by (3)
- (2) tends to be default output in the statistical software packages

Extremum estimation in some detail

Prerequisite: Uniform weak law of large numbers (UWLLN). This relates to certain random functions.

Definition (random function) Let (S, \mathcal{A}, P) be a probability space. A random function $f(\cdot, \boldsymbol{\theta})$ on a subset Θ of \mathbb{R}^m is a mapping

$$f(s, \boldsymbol{\theta}): S \times \Theta \rightarrow \mathbb{R}$$

such that \forall Borel set B in \mathbb{R} and $\forall \boldsymbol{\theta} \in \Theta$, the set $\{s \in S: f(s, \boldsymbol{\theta}) \in B\}$ is measurable \mathcal{A} , i.e. $P[\{s \in S: f(s, \boldsymbol{\theta}) \in B\}]$ is defined $B \in \mathcal{B}$ and $\forall \boldsymbol{\theta} \in \Theta$. So, random functions generally look like $g(\mathbf{x}, \boldsymbol{\theta})$, where \mathbf{x} is a random vector and $\boldsymbol{\theta}$ is a vector of constants.

Theorem 1 (UWLLN) Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be iid random k -vectors. Let $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$, where Θ is closed and bounded. Let $g(\mathbf{x}, \boldsymbol{\theta})$ be a Borel measurable function on $\mathbb{R}^k \times \Theta$ s.t. $g(\mathbf{x}, \boldsymbol{\theta})$ is continuous on $\Theta \forall \mathbf{x} \in \mathbb{R}^k$. Assume

$$E \left[\sup_{\boldsymbol{\theta} \in \Theta} |g(\mathbf{x}_1, \boldsymbol{\theta})| \right] < \infty$$

Then

$$\sup_{\boldsymbol{\theta} \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, \boldsymbol{\theta}) - E[g(\mathbf{x}_1, \boldsymbol{\theta})] \right| \xrightarrow{p} 0$$

i.e. $\frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, \boldsymbol{\theta})$ converges in probability to $E[g(\mathbf{x}_1, \boldsymbol{\theta})]$ uniformly on Θ .

Theorem 1 implies a result on the consistency of M-estimators based on iid observations.

Suppose the parameter of interest is

$$\boldsymbol{\theta}_0 \equiv \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} E[g(\mathbf{x}_1, \boldsymbol{\theta})]$$

where $g(\mathbf{x}, \boldsymbol{\theta})$ satisfies the conditions of Theorem 1.

$$\text{Let } \hat{\boldsymbol{\theta}}_n \equiv \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, \boldsymbol{\theta})$$

Theorem 2 (consistency of M-estimators) Assume

$$\mathbf{x}_1, \dots, \mathbf{x}_n, \quad g(\cdot, \cdot), \quad \Theta$$

are all as in the statement of Theorem 1. Assume that $\boldsymbol{\theta}_0$ is unique, i.e.

$$\forall \varepsilon > 0, \quad \exists \delta(\varepsilon) > 0, \quad \text{s.t. } E[g(\mathbf{x}_1, \boldsymbol{\theta}_0)] - \sup_{\boldsymbol{\theta}: \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| > \varepsilon} E[g(\mathbf{x}_1, \boldsymbol{\theta})] > \delta_\varepsilon$$

Then

$$\hat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}_0$$

Note: If Θ is compact and $E[g(\mathbf{x}_1, \boldsymbol{\theta})]$ is continuous on Θ then $E[g(\mathbf{x}_1, \boldsymbol{\theta})]$ is uniquely maximized on Θ .

Proof

Let

$$\begin{aligned}\widehat{Q}_n(\boldsymbol{\theta}) &\equiv \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i, \boldsymbol{\theta}) \\ \bar{Q}_n(\boldsymbol{\theta}) &\equiv E[g(\mathbf{x}_1, \boldsymbol{\theta})]\end{aligned}$$

Then $\widehat{\boldsymbol{\theta}}_n, \boldsymbol{\theta}_0 \in \Theta$ because $g(\mathbf{x}, \boldsymbol{\theta})$ is continuous in $\boldsymbol{\theta} \in \Theta$, which is compact. By the definition of $\boldsymbol{\theta}_0$,

$$\begin{aligned}0 &\leq \bar{Q}(\boldsymbol{\theta}_0) - \bar{Q}(\widehat{\boldsymbol{\theta}}_n) \\ &= \bar{Q}(\boldsymbol{\theta}_0) - \widehat{Q}_n(\boldsymbol{\theta}_0) + \widehat{Q}_n(\boldsymbol{\theta}_0) - \bar{Q}(\widehat{\boldsymbol{\theta}}_n) \\ &\leq \bar{Q}(\boldsymbol{\theta}_0) - \widehat{Q}_n(\boldsymbol{\theta}_0) + \widehat{Q}_n(\widehat{\boldsymbol{\theta}}_n) - \bar{Q}(\widehat{\boldsymbol{\theta}}_n) \\ &\leq 2 \sup_{\boldsymbol{\theta} \in \Theta} |\widehat{Q}_n(\boldsymbol{\theta}) - \bar{Q}(\boldsymbol{\theta})| \\ &\xrightarrow{p} 0, \quad \text{by Theorem 1,} \quad \dots (*)\end{aligned}$$

Therefore,

$$\bar{Q}(\widehat{\boldsymbol{\theta}}_n) \xrightarrow{p} \bar{Q}(\boldsymbol{\theta}_0)$$

But $\boldsymbol{\theta}_0$ is unique so $\forall \varepsilon > 0, \exists \delta(\varepsilon) > 0$ s.t. $\|\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\| > \varepsilon$ implies

$$\bar{Q}(\boldsymbol{\theta}_0) - \bar{Q}(\widehat{\boldsymbol{\theta}}_n) \geq \delta_\varepsilon$$

which implies

$$P[\|\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\| > \varepsilon] \leq P[\bar{Q}(\boldsymbol{\theta}_0) - \bar{Q}(\widehat{\boldsymbol{\theta}}_n) > \delta_\varepsilon] \rightarrow 0, \quad \text{by } (*), \quad \text{as } n \rightarrow \infty$$

$$\hat{\theta}_n \equiv \underset{\theta \in \Theta}{\operatorname{argmax}} \hat{Q}_n(\theta)$$

$$\hat{Q}_n(\theta) = \frac{1}{n} \sum_{j=1}^n g(\mathbf{X}_j, \theta)$$

Want to estimate θ_0 .

Assume sufficient conditions for consistency:

- $\Theta \subset \mathbb{R}^m$ is closed and bounded
- $g(\mathbf{x}, \theta)$ is Borel-measurable on $\mathbb{R}^k \times \Theta$ and $g(\mathbf{x}, \theta)$ is continuous in $\theta \forall \mathbf{x} \in \mathbb{R}^k$
- $E[\sup_{\theta \in \Theta} |g(\mathbf{X}_j, \theta)|] < \infty$
- $\mathbf{X}_1, \dots, \mathbf{X}_n$ are iid random k -vectors
- For $\bar{Q}(\theta) \equiv E[g(\mathbf{X}_1, \theta)]$ suppose that $\theta_0 \equiv \underset{\theta \in \Theta}{\operatorname{argmax}} \bar{Q}(\theta)$ is unique.

Add the following conditions that together with the conditions suffice for asymptotic Normality

- Θ is convex
- θ_0 is an interior point of Θ
- $\forall \mathbf{x} \in \mathbb{R}^k$, $g(\mathbf{x}, \theta)$ is twice continuously differentiable in $\theta \in \Theta$
- $E \left[\sup_{\theta \in \Theta} \left| \frac{\partial^2 g(\mathbf{X}_1, \theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right| \right] < \infty$ for any components (not necessarily distinct) $\theta_{i_1}, \theta_{i_2}$ of θ
- $\mathbf{A}_0 \equiv E \left[\frac{\partial^2 g(\mathbf{X}_1, \theta_0)}{\partial \theta_0 \partial \theta_0^T} \right]$ is nonsingular
- $\mathbf{B}_0 \equiv E \left[\frac{\partial g(\mathbf{X}_1, \theta_0)}{\partial \theta_0^T} \frac{\partial g(\mathbf{X}_1, \theta_0)}{\partial \theta_0} \right]$ is finite

Recall: $\hat{\theta}_n$ is asymptotically Normal iff there is an increasing sequence of positive numbers $\{a_n\}$ and a positive semidefinite matrix Σ s.t.

$$a_n(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(\mathbf{0}, \Sigma)$$

(this implies consistency of $\hat{\theta}_n$, since)

$$\hat{\theta}_n - \theta_0 = O_p(a_n^{-1}) = o_p(1)$$

Here we consider the case of extremum estimators that are root-n-consistent and asymptotically Normal, i.e.

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(\mathbf{0}, \Sigma)$$

Theorem

The eleven conditions given earlier imply

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} N_m(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

Proof

Consider only the case $m=1$. Recall that the first five conditions imply the consistency $\hat{\theta}_n \xrightarrow{p} \theta_0$. Since $\theta_0 \in \text{int}\Theta$, then $P[\hat{\theta}_n \in \text{int}\Theta] \rightarrow 1$. Therefore $P[\hat{Q}'_n(\hat{\theta}_n) = 0] \rightarrow 1$ (where we have applied differentiability of $\hat{Q}_n(\theta)$). Apply the two-times continuous differentiability of $\hat{Q}_n(\boldsymbol{\theta})$ and the mean-value theorem to deduce that $\exists \hat{\lambda}_n \in (0,1)$ s.t.

$$\sqrt{n}\hat{Q}'_n(\hat{\theta}_n) = \sqrt{n}\hat{Q}'_n(\theta_0) + \hat{Q}''_n(\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0))\sqrt{n}(\hat{\theta}_n - \theta_0)$$

Since Θ is convex, $P[\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0) \in \Theta] = 1$ and since $\hat{\theta}_n \xrightarrow{p} \theta_0$, we have $\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0) \xrightarrow{p} \theta_0$. By the UWLLN and the twice continuous differentiability of $\bar{Q}(\theta)$ (dominate convergence \rightarrow check, second-order partial derivatives?) we have

$$\sup_{\theta \in \Theta} |\hat{Q}''_n(\theta) - \bar{Q}''(\theta)| \xrightarrow{p} 0$$

Therefore $\hat{Q}''_n(\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0)) \xrightarrow{p} \bar{Q}''(\theta_0) < 0$ since θ_0 is the unique maximizer of $\bar{Q}(\theta)$. Therefore

$$\left(\hat{Q}''_n(\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0))\right)^{-1} \xrightarrow{p} (\bar{Q}''(\theta_0))^{-1} = \mathbf{A}_0^{-1}$$

Go back to the mean-value expansion, we have

$$\begin{aligned} \sqrt{n}(\hat{\theta}_n - \theta_0) &= -\left(\hat{Q}''_n(\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0))\right)^{-1} \sqrt{n}\hat{Q}'_n(\theta_0) + \left(\hat{Q}''_n(\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0))\right)^{-1} \sqrt{n}\hat{Q}'_n(\hat{\theta}_n) \\ &= -\left(\hat{Q}''_n(\theta_0 + \hat{\lambda}_n(\hat{\theta}_n - \theta_0))\right)^{-1} \sqrt{n}\hat{Q}'_n(\theta_0) + o_p(1) \end{aligned}$$

since $P[\hat{Q}'_n(\hat{\theta}_n) = 0] \rightarrow 1$. We have $\theta_0 \in \text{int}\Theta$, so

$$0 = \bar{Q}'(\theta_0) = E\left[\frac{dg(X_1, \theta_0)}{d\theta_0}\right]$$

In addition, we have

$$B_0 = \text{Var}\left[\frac{dg(X_1, \theta_0)}{d\theta_0}\right] \in (0, \infty), \quad \text{by assumption}$$

Therefore

$$\sqrt{n}\hat{Q}'_n(\theta_0) = \frac{1}{\sqrt{n}} \sum_{j=1}^n \frac{dg(X_j, \theta_0)}{d\theta_0} \xrightarrow{d} N(0, B_0)$$

$$\text{By Slutsky's theorem,} \quad \sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

In general, if an extremum estimator $\hat{\theta}_n$ satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_m(0, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

Need to estimate the asymptotic variance

$$\mathbf{V}[\hat{\theta}_n] = \frac{1}{n} \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1}$$

Let

$$\begin{aligned} \hat{\mathbf{A}}_n &\equiv \frac{1}{n} \sum_{j=1}^n \frac{\partial^2 g(\mathbf{X}_j, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n} \\ \hat{\mathbf{B}}_n &\equiv \frac{1}{n} \sum_{j=1}^n \frac{\partial g(\mathbf{X}_j, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n} \frac{\partial g(\mathbf{X}_j, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_n} \end{aligned}$$

Under the eleven conditions sufficient for the result

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_m(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

We have $\hat{\mathbf{A}}_n \xrightarrow{p} \mathbf{A}_0$

Impose the additional condition

$$E \left[\sup_{\boldsymbol{\theta} \in \Theta} \left\| \frac{\partial g(\mathbf{X}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\|^2 \right] < \infty$$

Then if the eleven conditions sufficient for asymptotic Normality hold, $\hat{\mathbf{B}}_n \xrightarrow{p} \mathbf{B}_0$ (exercise: check the applicability of the appropriate UWLLN). Therefore, under the twelve previous conditions

$$\hat{\mathbf{V}}_n[\hat{\theta}_n] \equiv \frac{1}{n} \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n \hat{\mathbf{A}}_n^{-1} \text{ satisfies } \|\hat{\mathbf{V}}_n[\hat{\theta}_n] - \mathbf{V}[\hat{\theta}_n]\| \xrightarrow{p} 0$$

So: If the twelve conditions hold, then it is straightforward to test hypothesis of the form

$$H_0: \mathbf{R}\boldsymbol{\theta}_0 = \mathbf{q}, \quad \text{vs } H_1: \mathbf{R}\boldsymbol{\theta}_0 \neq \mathbf{q}$$

where \mathbf{R} is an $(r \times m)$ -matrix of rank r , $\mathbf{q} \in \mathbb{R}^r$

To see this, under H_0 ,

$$\sqrt{n}(\mathbf{R}\hat{\boldsymbol{\theta}}_n - \mathbf{q}) \xrightarrow{d} N_r(\mathbf{0}, \mathbf{R}\mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1} \mathbf{R}^\top)$$

(also if \mathbf{B}_0 is nonsingular, then $\mathbf{R}\mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1} \mathbf{R}^\top$ is invertible, otherwise can consider a generalized inverse).

By Slutsky's theorem:

$$W_n = n(\mathbf{R}\hat{\boldsymbol{\theta}}_n - \mathbf{q})^\top (\mathbf{R}\hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n \hat{\mathbf{A}}_n^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\hat{\boldsymbol{\theta}}_n - \mathbf{q}) \xrightarrow{d} \chi_r^2, \quad \text{under } H_0$$

This allows one to formulate a test of (approximate) level α when n is large.

We have $P[W_n > B_\alpha] \rightarrow \alpha$ under H_0 as $n \rightarrow \infty$, where $B_\alpha = (1-\alpha)$ -quantile of a χ_r^2 .

So: If rejection of $H_0: \mathbf{R}\boldsymbol{\theta}_0 = \mathbf{q}$ occurs iff $W_n > B_\alpha$ then the probability of type-I error is approximately α when n is large.

To see that this “Wald-type” test has power, suppose H_0 is false, then under our twelve conditions

$$\frac{W_n}{n} \xrightarrow{p} (\mathbf{R}\boldsymbol{\theta}_0 - \mathbf{q})^\top (\mathbf{R}\mathbf{A}_0^{-1}\mathbf{B}_0\mathbf{A}_0^{-1}\mathbf{R}^\top)^{-1} (\mathbf{R}\boldsymbol{\theta}_0 - \mathbf{q}) > 0$$

e.g. Poisson MLE: Suppose we observe $(\mathbf{X}_1^\top, Y_1)^\top, \dots, (\mathbf{X}_n^\top, Y_n)^\top$ iid. Assume $E[Y_i|\mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ w.p. 1. Let

$$\widehat{Q}_n(\boldsymbol{\beta}) \equiv \frac{1}{n} \sum_{i=1}^n \left(-e^{\mathbf{X}_i^\top \boldsymbol{\beta}} + Y_i \mathbf{X}_i^\top \boldsymbol{\beta} - \log Y_i! \right), \quad \text{i. e. the Poisson likelihood}$$

Let

$$\bar{Q}(\boldsymbol{\beta}) \equiv E[\widehat{Q}_n(\boldsymbol{\beta})]$$

By Kolmogorov’s LLN

$$\bar{Q}(\boldsymbol{\beta}) = \text{plim } \widehat{Q}_n(\boldsymbol{\beta}) = -E[e^{\mathbf{X}_1^\top \boldsymbol{\beta}}] + E[Y_1 \mathbf{X}_1^\top \boldsymbol{\beta}] - E[\log Y_1!]$$

provided that each of these expectations exists. Since $E[Y_i|\mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ w.p. 1, so

$$\bar{Q}(\boldsymbol{\beta}) = -E[e^{\mathbf{X}_1^\top \boldsymbol{\beta}}] + E[e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0} \mathbf{X}_1^\top \boldsymbol{\beta}] - E[\log Y_1!]$$

Therefore, if dominated convergence is applicable (say $\boldsymbol{\beta} \in B$, where B is compact and $E[\|\mathbf{X}_1\|^2] < \infty$) then

$$\frac{\partial}{\partial \boldsymbol{\beta}^\top} \bar{Q}(\boldsymbol{\beta}) = -E[\mathbf{X}_1 e^{\mathbf{X}_1^\top \boldsymbol{\beta}}] + E[\mathbf{X}_1 e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0}], \quad \text{i. e. } \frac{\partial}{\partial \boldsymbol{\beta}^\top} \bar{Q}(\boldsymbol{\beta}) \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} = \mathbf{0}$$

In addition,

$$\frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \bar{Q}(\boldsymbol{\beta}) = -E[\mathbf{X}_1 \mathbf{X}_1^\top e^{\mathbf{X}_1^\top \boldsymbol{\beta}}]$$

which is negative definite if there is no perfect multicollinearity.

Therefore $\bar{Q}(\boldsymbol{\beta})$ is globally maximized at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$.

Consistency of the Poisson MLE follows from what we’ve discussed.

What about asymptotic Normality? Take a Taylor expansion to first-order of the FOCs to get

$$\frac{1}{n} \sum_{i=1}^n \left(Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}} \right) \mathbf{X}_i \Big|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_n} = \mathbf{0}$$

to get

$$\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) = - \left(-\frac{1}{n} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\beta}}_n} \mathbf{X}_i \mathbf{X}_i^\top \right)^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0} \right) \mathbf{X}_i$$

where $\bar{\boldsymbol{\beta}}_n$ is between $\widehat{\boldsymbol{\beta}}_n$ and $\boldsymbol{\beta}_0$. We have $\widehat{\boldsymbol{\beta}}_n \xrightarrow{p} \boldsymbol{\beta}_0$, so $\bar{\boldsymbol{\beta}}_n \xrightarrow{p} \boldsymbol{\beta}_0$, so if $E[\|\mathbf{X}_1\|^4] < \infty$ (why?) then

$$-\frac{1}{n} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\beta}}_n} \mathbf{X}_i \mathbf{X}_i^\top \xrightarrow{p} -E[e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0} \mathbf{X}_1 \mathbf{X}_1^\top] \equiv \mathbf{A}_0, \quad \dots (*)$$

Next, suppose \mathbf{X}_i is a scalar, i.e. $\mathbf{X}_i = X_i$. Then

$$E[(Y_1 - e^{X_1 \beta_0}) X_1] = 0$$

and

$$\begin{aligned} \text{Var}[(Y_1 - e^{X_1 \beta_0}) X_1] &= E[\text{Var}[(Y_1 - e^{X_1 \beta_0}) X_1 | X_1]] + \text{Var}[E[(Y_1 - e^{X_1 \beta_0}) X_1 | X_1]] \\ &= E[\text{Var}[Y_1 | X_1] X_1^2] < \infty \end{aligned}$$

if $E[X_1^2] < \infty$ and $\text{Var}[Y_1 | X_1] < \infty$. By Lindberg–Lévy CLT,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i - e^{X_i \beta_0}) X_i \xrightarrow{d} N(0, E[\text{Var}[Y_1 | X_1] X_1^2]), \quad \dots (**)$$

Can extend (*) to the case where \mathbf{X}_1 is vector-valued by Cramér–Wold device.

In particular,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \left(Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0} \right) \mathbf{X}_i \xrightarrow{d} N_d(\mathbf{0}, E[\text{Var}[Y_1 | \mathbf{X}_1] \mathbf{X}_1 \mathbf{X}_1^\top])$$

Let

$$\mathbf{B}_0 \equiv E[\text{Var}[Y_1 | \mathbf{X}_1] \mathbf{X}_1 \mathbf{X}_1^\top]$$

Combining results, we get

$$\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

by Slutsky's theorem.

Poisson MLE

Observe $[\mathbf{X}_i^\top Y_i]$ ($i=1, \dots, n$) iid

$$\begin{aligned}\widehat{Q}_n(\boldsymbol{\beta}) &\equiv \frac{1}{n} \sum_{i=1}^n \left(-e^{\mathbf{X}_i^\top \boldsymbol{\beta}} + Y_i \mathbf{X}_i^\top \boldsymbol{\beta} - \log Y_i! \right) \\ \bar{Q}(\boldsymbol{\beta}) &\equiv \text{plim } \widehat{Q}_n(\boldsymbol{\beta}) \\ &= E[\widehat{Q}_n(\boldsymbol{\beta})] \\ &= -E \left[e^{\mathbf{X}_i^\top \boldsymbol{\beta}} \right] + E[Y_i \mathbf{X}_i^\top \boldsymbol{\beta}] - E[\log Y_i!]\end{aligned}$$

Assume $E[Y_1 | \mathbf{X}_1] = e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0}$ w.p. 1. We showed that

$$\left. \frac{\partial}{\partial \boldsymbol{\beta}^\top} \bar{Q}(\boldsymbol{\beta}) \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} = \mathbf{0}$$

(assuming that $E[Y_1 | \mathbf{X}_1] = e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0}$ and dominated convergence applies)

Also

$$\left. \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \bar{Q}(\boldsymbol{\beta}) \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} = -E \left[e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0} \mathbf{X}_1 \mathbf{X}_1^\top \right]$$

i.e. negative definite if no perfect multicollinearity.

Therefore under appropriate conditions $\bar{Q}(\boldsymbol{\beta})$ is globally maximized at $\boldsymbol{\beta}=\boldsymbol{\beta}_0$ and consistency of the Poisson MLE follows.

What about asymptotic Normality?

Take a Taylor expansion to first order of the left hand side of the system

$$\left. \frac{1}{n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}}) \mathbf{X}_i \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_n} = \mathbf{0}$$

to get

$$\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) = - \left(-\frac{1}{n} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\beta}}_n} \mathbf{X}_i \mathbf{X}_i^\top \right)^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}) \mathbf{X}_i$$

where $\bar{\boldsymbol{\beta}}_n$ is between $\widehat{\boldsymbol{\beta}}_n$ and $\boldsymbol{\beta}_0$.

Since $\widehat{\boldsymbol{\beta}}_n \xrightarrow{p} \boldsymbol{\beta}_0$, $\bar{\boldsymbol{\beta}}_n \xrightarrow{p} \boldsymbol{\beta}_0$ and therefore if e.g. $E[\|\mathbf{X}_1\|^4] < \infty$ we get

$$-\frac{1}{n} \sum_{i=1}^n e^{\mathbf{X}_i^\top \bar{\boldsymbol{\beta}}_n} \mathbf{X}_i \mathbf{X}_i^\top \xrightarrow{p} -E \left[e^{\mathbf{X}_1^\top \boldsymbol{\beta}_0} \mathbf{X}_1 \mathbf{X}_1^\top \right] \equiv \mathbf{A}_0, \quad \left(\begin{array}{l} \text{exercise:} \\ \text{supply details} \end{array} \right)$$

Next suppose \mathbf{X}_i is scalar, then

$$E[(Y_1 - e^{X_1\beta_0})X_1] = 0$$

and

$$\text{Var}[(Y_1 - e^{X_1\beta_0})X_1] = E[\text{Var}[Y_1|X_1]X_1^2], \quad \left(\begin{array}{l} \text{exercise:} \\ \text{supply details} \end{array} \right)$$

Assume $P[\text{Var}[Y_1|X_1] < \infty] = 1$, $E[X_1^2] < \infty$. Then by the Lindberg–Lévy CLT

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n (Y_i - e^{X_i\beta_0})X_i \xrightarrow{d} N(0, E[\text{Var}[Y_1|X_1]X_1^2])$$

The case with \mathbf{X}_i not necessarily scalar follows the Cramér–Wold device (exercise: review).

Let $\mathbf{B}_0 \equiv E[\text{Var}[Y_1|\mathbf{X}_1]\mathbf{X}_1\mathbf{X}_1^\top]$

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{A}_0^{-1}\mathbf{B}_0\mathbf{A}_0^{-1}), \quad \dots (*)$$

Important: Convergence (*) holds even if $Y_i|\mathbf{X}_i$ is not $\text{Poisson}(e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0})$ all that is required for consistency that $E[Y_i|\mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ (along with other conditions like $E[\|\mathbf{X}_1\|^4] < \infty$). Asymptotic Normality also requires that $\text{Var}[Y_i|\mathbf{X}_i] < \infty$ w.p. 1. On the other hand, if $\text{Var}[Y_i|\mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ (for example, if $Y_i|\mathbf{X}_i \sim \text{Poisson}(e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0})$ then $\mathbf{A}_0 = -\mathbf{B}_0$ and the convergence in (*) simplifies to)

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{d} N(\mathbf{0}, -\mathbf{A}_0^{-1})$$

In general, the Poisson MLE $\hat{\boldsymbol{\beta}}_n$ is really a quasi-MLE, i.e. don't need to assume the conditional distribution of $Y_i|\mathbf{X}_i$ in order to estimate parameter of $E[Y_i|\mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$.

But: Assuming correctness of the nonlinear regression $E[Y_i|\mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ is the asymptotic variance in (*) the best possible?

(In other words, is it possible for another M-estimator of \mathbf{B}_0 to be \sqrt{n} -consistent and asymptotically Normal, but generate smaller standard errors?)

For example, nonlinear least squares, i.e.

$$\hat{\boldsymbol{\beta}}_{\text{NLLS}} = \underset{\boldsymbol{\beta}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}})^2$$

The answer to this question involves the theory of efficient GMM.

On the other hand, suppose $Y_i|\mathbf{X}_i \sim \text{Poisson}(e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0})$ holds for each i .

In this case, the MLE $\hat{\boldsymbol{\beta}}_n$ is asymptotically efficient its asymptotic variance $-\frac{1}{n}\mathbf{A}_0^{-1}$ dominates the asymptotic variance of any other “regular” estimator of \mathbf{B}_0 .

(we call $\tilde{\boldsymbol{\beta}}_n$ is regular estimator of $\boldsymbol{\beta}_0$ if it is \sqrt{n} -consistent and asymptotically Normal with $\sqrt{n}(\tilde{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{V}(\boldsymbol{\beta}))$ for all $\boldsymbol{\beta}$ in an open neighborhood of $\boldsymbol{\beta}_0$, i.e. \sqrt{n} -consistency and asymptotic Normality of regular estimators are maintained for all $\boldsymbol{\beta}$ near $\boldsymbol{\beta}_0$)

Theory of Maximum Likelihood

e.g. Let (S, \mathcal{A}, P) be a probability space

Let Z_1, \dots, Z_n be random variables defined on (S, \mathcal{A}, P)

Where Z_1, \dots, Z_n are iid Uniform $(0, \theta_0)$ for some $\theta_0 > 0$

The common density of the Z_i s

$$f(z, \theta_0) \equiv \frac{1}{\theta_0} 1\{0 \leq z \leq \theta_0\}$$

So the likelihood function is

$$\hat{L}_n(\theta) \equiv \theta^{-1} \prod_{j=1}^n 1\{0 \leq z_j \leq \theta\}, \quad \left(= \prod_{j=1}^n f(z_j, \theta) \right)$$

e.g. Let $\mathbf{Z}_j \equiv [\mathbf{X}_j^\top \ Y_j]^\top$ ($j=1, \dots, n$) be independent with

$$Y_j = \alpha_0 + \mathbf{X}_j^\top \boldsymbol{\beta}_0 + u_j, \quad \text{where } u_j | \mathbf{X}_j \sim N(0, \sigma_0^2)$$

(i.e. Normal linear model)

The conditional density of Y_j given \mathbf{X}_j is

$$f(y | \mathbf{X}_j; \boldsymbol{\theta}_0) = \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(y - \alpha_0 - \mathbf{X}_j^\top \boldsymbol{\beta}_0)^2}{2\sigma_0^2}}$$

Assume that the \mathbf{X}_j s are iid with absolutely continuous distribution having density $g(\cdot)$. Then the likelihood function

$$\begin{aligned} \hat{L}_n(\boldsymbol{\theta}) &= \prod_{j=1}^n f(Y_j | \mathbf{X}_j; \boldsymbol{\theta}) g(\mathbf{X}_j) \\ &= \frac{1}{(\sqrt{2\pi}\sigma)^n} e^{-\frac{\sum_{j=1}^n (Y_j - \alpha - \mathbf{X}_j^\top \boldsymbol{\beta})^2}{2\sigma^2}} \prod_{j=1}^n g(\mathbf{X}_j) \\ \text{where } \boldsymbol{\theta} &= [\alpha \ \boldsymbol{\beta}^\top \ \sigma^2]^\top \end{aligned}$$

Assuming that $g(\cdot)$ does not depend on $\boldsymbol{\theta}_0$ can work with the conditional likelihood function

$$\hat{L}_n^c(\boldsymbol{\theta}) = \prod_{j=1}^n f(Y_j | \mathbf{X}_j; \boldsymbol{\theta})$$

In general, writing down a likelihood function is nontrivial.

- The data may not be iid.
- The joint distribution of the data may not be absolutely continuous or discrete.

Better to define a likelihood function implicitly.

Definition 1 (likelihood function) A sequence $\{\hat{L}_n(\boldsymbol{\theta}): n=1,2,\dots\}$ of a nonnegative random functions in $\boldsymbol{\theta} \in \Theta$ is a sequence of likelihood functions iff

L.1. \exists an increasing sequence $\{\mathcal{A}_n: n = 1,2, \dots\}$ of σ -algebras such that $\forall \boldsymbol{\theta} \in \Theta$ and $n \geq 1$, $\hat{L}_n(\boldsymbol{\theta})$ is measurable \mathcal{A}_n

L.2. $\exists \boldsymbol{\theta}_0 \in \Theta$ s.t. $\forall \boldsymbol{\theta} \in \Theta$

$$P \left[E \left[\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_n(\boldsymbol{\theta})} \middle| \mathcal{A}_0 \right] \leq 1 \right] = 1$$

and $\forall n \geq 2$

$$P \left[E \left[\frac{\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_{n-1}(\boldsymbol{\theta})}}{\frac{\hat{L}_n(\boldsymbol{\theta}_0)}{\hat{L}_{n-1}(\boldsymbol{\theta}_0)}} \middle| \mathcal{A}_{n-1} \right] \leq 1 \right] = 1$$

L.3. $\forall \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2$,

$$P[\hat{L}_1(\boldsymbol{\theta}_1) = \hat{L}_1(\boldsymbol{\theta}_2) | \mathcal{A}_0] < 1$$

and $n \geq 2$

$$P \left[\frac{\hat{L}_n(\boldsymbol{\theta}_1)}{\hat{L}_{n-1}(\boldsymbol{\theta}_1)} = \frac{\hat{L}_n(\boldsymbol{\theta}_2)}{\hat{L}_{n-1}(\boldsymbol{\theta}_2)} \middle| \mathcal{A}_{n-1} \right] < 1$$

(L3 rules out functions $\hat{L}_n(\boldsymbol{\theta})$ that are w.p. 1 on Θ)

L2, L3 guarantee that the point $\boldsymbol{\theta}_0 \in \Theta$ (defined in L2) is unique.

Theorem 1 $\forall \boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}_0\}$ and $\forall n \geq 1$,

$$E \left[\log \left(\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_n(\boldsymbol{\theta}_0)} \right) \right] < 0$$

Proof

Let $n=1$, we have

$$\log \left(\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} \right) < \frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} - 1, \quad \text{by the property of } \log(\cdot)$$

if $\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} \neq 1$

Let

$$Y(\boldsymbol{\theta}) \equiv \frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} - \log \left(\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} \right) - 1, \quad X(\boldsymbol{\theta}) \equiv \frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)}$$

Then

$$Y(\boldsymbol{\theta}) \geq 0 \text{ with } Y(\boldsymbol{\theta}) > 0 \text{ iff } X(\boldsymbol{\theta}) \neq 1$$

Suppose

$$P[E[Y(\boldsymbol{\theta})|\mathcal{A}_0] = 0] = 1$$

Then

$$P[Y(\boldsymbol{\theta}) = 0|\mathcal{A}_0] = 1 \text{ a.s.}$$

because

$$Y(\boldsymbol{\theta}) \geq 0, \quad > 0 \text{ in this case,} \quad P[X(\boldsymbol{\theta}) = 1|\mathcal{A}_0] = 1 \text{ a.s.}$$

By L3 it cannot be true that $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ so

$$P \left[E \left[\log \left(\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} \right) \middle| \mathcal{A}_0 \right] < 0 \right] = 1 \text{ iff } \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$$

So if $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$, then

$$E \left[\log \left(\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} \right) \right] < 0$$

A similar argument shows that if $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$, then

$$E \left[\log \left(\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_{n-1}(\boldsymbol{\theta})} \right) - \log \left(\frac{\hat{L}_n(\boldsymbol{\theta}_0)}{\hat{L}_{n-1}(\boldsymbol{\theta}_0)} \right) \right] < 0$$

Definition (likelihood function) A sequence $\{\hat{L}_n(\boldsymbol{\theta}): n=1,2,\dots\}$ of non-negative random functions on a parameter space Θ is a sequence of likelihood functions iff

L1 \exists increasing sequence $\{\mathcal{A}_n: n=0,1,2,\dots\}$ of σ -algebras such that $\forall \boldsymbol{\theta} \in \Theta$ and $\forall n \geq 1$ $\hat{L}_n(\boldsymbol{\theta})$ is measurable \mathcal{A}_n .

L2 $\boldsymbol{\theta}_0 \in \Theta$ such that

$$P \left[E \left[\frac{\hat{L}_1(\boldsymbol{\theta})}{\hat{L}_1(\boldsymbol{\theta}_0)} \middle| \mathcal{A}_0 \right] \leq 1 \right] = 1$$

and $\forall n \geq 2$

$$P \left[E \left[\frac{\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_{n-1}(\boldsymbol{\theta})}}{\frac{\hat{L}_n(\boldsymbol{\theta}_0)}{\hat{L}_{n-1}(\boldsymbol{\theta}_0)}} \middle| \mathcal{A}_{n-1} \right] \leq 1 \right] = 1$$

L3 $\forall \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2, P[\hat{L}_1(\boldsymbol{\theta}_1) = \hat{L}_1(\boldsymbol{\theta}_2) | \mathcal{A}_0] < 1$, and $\forall n \geq 2$

$$P \left[\frac{\hat{L}_n(\boldsymbol{\theta}_1)}{\hat{L}_{n-1}(\boldsymbol{\theta}_1)} = \frac{\hat{L}_n(\boldsymbol{\theta}_2)}{\hat{L}_{n-1}(\boldsymbol{\theta}_2)} \middle| \mathcal{A}_{n-1} \right] < 1$$

(L3 rules out functions $\hat{L}_n(\boldsymbol{\theta})$ that are almost surely constant in $\boldsymbol{\theta}$ on Θ)

L3 implies that the point $\boldsymbol{\theta}_0$ defined in L2 is unique.

Theorem 1 $\forall \boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}_0\}$ and $\forall n \geq 1$

$$E \left[\log \left(\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_n(\boldsymbol{\theta}_0)} \right) \right] < 0$$

Also, if the data are iid from an absolutely continuous distribution with density $f(\cdot; \boldsymbol{\theta})$ the assumption that

$$\text{supp } f(\cdot; \boldsymbol{\theta}) \equiv \{\mathbf{z} \in \mathbb{R}^k: f(\mathbf{z}; \boldsymbol{\theta}) > 0\}$$

is the same for all $\boldsymbol{\theta} \in \Theta$ implies that

$$E \left[\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_n(\boldsymbol{\theta}_0)} \right] = 1$$

$\forall \boldsymbol{\theta} \in \Theta$ and all $n \geq 1$; this implication also goes in the other direction.

Definition 2 (invariant support) A sequence $\{\hat{L}_n(\theta); n=1,2,\dots\}$ of likelihood functions has invariant support iff $\forall \theta \in \Theta$

$$P \left[E \left[\frac{\hat{L}_1(\theta)}{\hat{L}_1(\theta_0)} \middle| \mathcal{A}_0 \right] = 1 \right] = 1$$

and $\forall n \geq 2$

$$P \left[E \left[\frac{\frac{\hat{L}_n(\theta)}{\hat{L}_{n-1}(\theta)}}{\frac{\hat{L}_n(\theta_0)}{\hat{L}_{n-1}(\theta_0)}} \middle| \mathcal{A}_{n-1} \right] = 1 \right] = 1$$

Note: Most likelihood functions in Econometrics have invariant support.

In general, L1-L3 allow one to treat MLE as special case of extremum estimation with

$$\theta_0 = \underset{\theta \in \Theta}{\operatorname{argmax}} E[\log \hat{L}_n(\theta)]$$

(recall: need L2, L3 for this)

This may require sometimes one to choose Θ sufficiently small.

e.g. Suppose $Z_j = \cos(X_j + \theta_0)$ where X_1, X_2, \dots are iid with an absolutely continuous distribution with density f .

Then the density of Z_j has the form $f(z; \theta_0) = f(z; \theta_0 + 2\pi s)$ for any $s \in \mathbb{Z} \equiv$ set of integers.

In this case, need Θ to be small, say $\Theta = [0, 2\pi]$ in order to make θ_0 unique.

In addition: MLE cannot be an extremum estimator if $E[\log \hat{L}_n(\theta)]$ does not satisfy first- and second-order conditions for the existence of a global maximum at $\theta = \theta_0$ when θ_0 is in the interior of Θ .

However, suppose θ_0 is on the boundary of Θ .

e.g. Z_1, \dots, Z_n are iid Uniform $(0, \theta_0)$ for some $\theta_0 > 0$. Then we know

$$\hat{L}_n(\theta) = \frac{1}{\theta^n} \prod_{j=1}^n 1\{Z_j \in [0, \theta]\}$$

We have

$$E[\log \hat{L}_n(\theta)] = \begin{cases} -\infty, & \theta < \theta_0 \\ -n \log \theta, & \theta \geq \theta_0 \end{cases}$$

The left derivative of $E[\log \hat{L}_n(\theta)]$ at $\theta = \theta_0$ is

$$\lim_{\delta \uparrow 0} \frac{1}{\delta} (E[\log \hat{L}_n(\theta_0)] - E[\log \hat{L}_n(\theta_0 - \delta)]) = \infty$$

The right derivative of $E[\log \hat{L}_n(\theta)]$ at $\theta = \theta_0$ is

$$\lim_{\delta \downarrow 0} \frac{1}{\delta} (E[\log \hat{L}_n(\theta_0 + \delta)] - E[\log \hat{L}_n(\theta_0)]) = -\frac{n}{\theta_0}$$

So $E[\log \hat{L}_n(\theta)]$ is not differentiable at θ_0 .

Therefore the MLE $\hat{\theta}_n = Z_{(n)}$ of θ_0 does not behave as a standard extremum estimator; cannot derive its asymptotic properties from first and second order conditions.

Henceforth: Only consider MLEs for problems where θ_0 is defined in first- and second-order conditions.

Assumption 1 Θ is convex; $\theta_0 \in \text{interior}\Theta$. Also $\hat{L}_n(\theta)$ is (w.p. 1) twice continuously differentiable in an open neighborhood Θ_0 of θ_0 .

If $\theta_0 \in \mathbb{R}^d$ then for each $i_1, i_2 \in \{1, \dots, d\}$ assume

$$E \left[\sup_{\theta \in \Theta_0} \left| \frac{\partial^2 \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right| \right] < \infty, \quad \dots (22)$$

and

$$E \left[\sup_{\theta \in \Theta_0} \left| \frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right| \right] < \infty, \quad \dots (23)$$

Theorem 2 (score; information identities)

$$E \left[\left. \frac{\partial \log \hat{L}_n(\theta)}{\partial \theta^\top} \right|_{\theta=\theta_0} \right] = \mathbf{0}, \quad \dots (\text{score identity})$$

$$E \left[\left. \frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta \partial \theta^\top} \right|_{\theta=\theta_0} \right] = -\text{Var} \left[\left. \frac{\partial \log \hat{L}_n(\theta)}{\partial \theta^\top} \right|_{\theta=\theta_0} \right], \quad \dots (\text{information identity})$$

(Note: Without the score identity, the MLE is inconsistent. If the score identity holds, but the information identity fails, the MLE is consistent but asymptotically inefficient)

Proof Case $d=1$ only. Also assume that $\mathbf{Z} \equiv [\mathbf{z}_1 \ \dots \ \mathbf{z}_n]^\top$ is a random sample from an absolutely continuous distribution with density $f(\mathbf{z}; \theta_0)$. Then

$$E \left[\frac{1}{n} \log \hat{L}_n(\theta) \right] = \frac{1}{n} \sum_{i=1}^n E[\log f(\mathbf{z}_i; \theta)] = \int (\log f(\mathbf{z}; \theta)) f(\mathbf{z}; \theta_0) d\mathbf{z}$$

By Taylor's theorem, $\forall \theta \in \Theta_0$ and $\forall \delta \neq 0$ with $\theta + \delta \in \Theta_0$, $\exists \lambda(\mathbf{z}, \delta) \in [0, 1]$ such that

$$\log f(\mathbf{z}; \theta + \delta) = \log f(\mathbf{z}; \theta) + \delta \frac{d}{d\theta} \log f(\mathbf{z}; \theta) + \frac{\delta^2}{2} \frac{d^2}{d\tilde{\theta}^2} \log f(\mathbf{z}; \tilde{\theta}) \Big|_{\tilde{\theta}=\theta+\lambda(\mathbf{z};\delta)\delta}$$

Since Θ is convex $\theta + \lambda(\mathbf{z}, \delta) \delta \in \Theta$. By (23) and the dominated convergence theorem

$$\frac{d}{d\theta} \int (\log f(\mathbf{z}; \theta)) f(\mathbf{z}; \theta_0) d\mathbf{z} = \int \frac{d \log f(\mathbf{z}; \theta)}{d\theta} f(\mathbf{z}; \theta_0) d\mathbf{z}, \quad \dots (26)$$

By condition (22) and the dominated convergence theorem, we have

$$\int \frac{df(\mathbf{z}; \theta)}{d\theta} d\mathbf{z} = \frac{d}{d\theta} \int f(\mathbf{z}; \theta) d\mathbf{z} = \frac{d}{d\theta} 1 = 0, \quad \dots (27)$$

Also,

$$\begin{aligned} \int \left. \frac{d \log f(\mathbf{z}; \theta)}{d\theta} \right|_{\theta=\theta_0} f(\mathbf{z}; \theta_0) d\mathbf{z} &= \int \frac{1}{f(\mathbf{z}; \theta)} \left. \frac{df(\mathbf{z}; \theta)}{d\theta} \right|_{\theta=\theta_0} f(\mathbf{z}; \theta_0) d\mathbf{z} \\ &= \frac{d}{d\theta} \int f(\mathbf{z}; \theta_0) d\mathbf{z}, \quad \dots (26) \end{aligned}$$

Combine (24)–(28) to get

$$E \left[\left. \frac{d}{d\theta} \frac{1}{n} \log \hat{L}_n(\theta) \right|_{\theta=\theta_0} \right] = E \left[\left. \frac{d \log f(\mathbf{z}_1; \theta)}{d\theta} \right|_{\theta=\theta_0} \right] = \int \left(\left. \frac{d}{d\theta} \log f(\mathbf{z}; \theta) \right|_{\theta=\theta_0} \right) f(\mathbf{z}; \theta_0) d\mathbf{z} = 0$$

i.e. the score identity.

Next: Consider the information identity. By (22), (23) and the dominated convergence theorem we have

$$\frac{d^2}{d\theta^2} \int (\log f(\mathbf{z}; \theta)) f(\mathbf{z}; \theta_0) d\mathbf{z} = \int \left(\frac{d^2}{d\theta^2} \log f(\mathbf{z}; \theta) \right) f(\mathbf{z}; \theta_0) d\mathbf{z}, \quad \dots (29)$$

and

$$\int \frac{d^2 f(\mathbf{z}; \theta)}{d\theta^2} d\mathbf{z} = \frac{d^2}{d\theta^2} \int f(\mathbf{z}; \theta) d\mathbf{z} = 0, \quad \dots (30)$$

(29), (30) imply the information identity because

$$\begin{aligned} \int \left(\left. \frac{d^2}{d\theta^2} \log f(\mathbf{z}; \theta) \right|_{\theta=\theta_0} \right) f(\mathbf{z}; \theta_0) d\mathbf{z} &= \int \left(\left. \frac{\frac{d^2}{d\theta^2} f(\mathbf{z}; \theta)}{f(\mathbf{z}; \theta)} \right|_{\theta=\theta_0} \right) f(\mathbf{z}; \theta_0) d\mathbf{z} \\ &\quad - \int \left(\left. \left(\frac{\frac{d}{d\theta} f(\mathbf{z}; \theta)}{f(\mathbf{z}; \theta)} \right)^2 \right|_{\theta=\theta_0} \right) f(\mathbf{z}; \theta_0) d\mathbf{z} \\ &= \int \left. \frac{d^2}{d\theta^2} f(\mathbf{z}; \theta) \right|_{\theta=\theta_0} d\mathbf{z} - \int \left(\left. \left(\frac{d}{d\theta} \log f(\mathbf{z}; \theta) \right)^2 \right|_{\theta=\theta_0} \right) f(\mathbf{z}; \theta_0) d\mathbf{z} \\ &= - \int \left(\left. \frac{d}{d\theta} \log f(\mathbf{z}; \theta) \right|_{\theta=\theta_0} \right)^2 f(\mathbf{z}; \theta_0) d\mathbf{z} \\ &= - \text{Var} \left[\left. \frac{d}{d\theta} \log f(\mathbf{z}_1; \theta) \right|_{\theta=\theta_0} \right] \text{ by the score identity} \end{aligned}$$

i.e. the information identity.

In general, the (d×d) matrix

$$\mathbf{H} \equiv \text{Var} \left[\frac{\partial}{\partial \boldsymbol{\theta}^\top} \log \hat{L}_n(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right]$$

is the Fischer information (matrix) for $\boldsymbol{\theta}_0$.

(Recall: If θ is scalar and the observations are iid, then H^{-1} is the Cramér–Rao lower bound for the variance of an unbiased estimator of θ_0)

Consistency and Asymptotic Normality of MLEs

- Consider MLE as a special case of extremum estimation
- Need to assume that

$$\boldsymbol{\theta}_0 = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} E[\log \hat{L}_n(\boldsymbol{\theta})]$$

can be characterized in terms of first- and second-order conditions.

Also need the following.

Assumption 2 As $n \rightarrow \infty$

$$\sup_{\boldsymbol{\theta} \in \Theta} \left| \frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta}) - E[\log \hat{L}_n(\boldsymbol{\theta})] \right| \xrightarrow{p} 0$$

and

$$\sup_{\boldsymbol{\theta} \in \Theta} \left| \frac{1}{n} E \left[\log \left(\frac{\hat{L}_n(\boldsymbol{\theta})}{\hat{L}_n(\boldsymbol{\theta}_0)} \right) \right] - \ell(\boldsymbol{\theta}; \boldsymbol{\theta}_0) \right| \rightarrow 0$$

for some continuous function $\ell(\boldsymbol{\theta}; \boldsymbol{\theta}_0)$ of $\boldsymbol{\theta}$ s.t. $\forall \delta > 0$

$$\sup_{\boldsymbol{\theta} \in \Theta, \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \geq \delta} \ell(\boldsymbol{\theta}; \boldsymbol{\theta}_0) < 0, \quad \dots (*)$$

Theorem 3 (consistency) Under assumptions 1, 2

$$\hat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}$$

Note: Assumption 2 can only be verified in the context of specific examples.

- Can consider conditions that suffice for a ULLN: e.g. Θ compact, $\log \hat{L}_n(\boldsymbol{\theta})$ to be continuous and uniformly bounded on Θ w.p. 1
- Condition (*) in Assumption 2 will hold if Θ is compact and $\ell(\boldsymbol{\theta}; \boldsymbol{\theta}_0)$ is continuous on Θ and $\boldsymbol{\theta}_0$ is unique ($\boldsymbol{\theta}_0$ is unique by Theorem 1; Theorem 1 holds by L2, L3)

Recall certain assumptions that allow one to treat MLE as “just another” M-estimator

Assumption 1

- The parameter space Θ is convex
- $\theta_0 \in \text{int}\Theta$
- The likelihood function $\hat{L}_n(\theta)$ is (w.p. 1) twice continuously differentiable in an open neighborhood Θ_0 of θ_0
- If $\theta_0 \in \mathbb{R}^d$, then for each pair $i_1, i_2 \in \{1, \dots, d\}$

$$E \left[\sup_{\theta \in \Theta_0} \left| \frac{\partial^2 \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right| \right] < \infty, \quad \text{and } E \left[\sup_{\theta \in \Theta_0} \left| \frac{\partial^2 \ln \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right| \right] < \infty$$

(Recall that Assumption 1 implies the score and information identities: score identity \rightarrow consistency, information identity \rightarrow asymptotic efficiency)

Assumption 2 As $n \rightarrow \infty$

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \log \hat{L}_n(\theta) - E \left[\frac{1}{n} \log \hat{L}_n(\theta) \right] \right| \xrightarrow{p} 0, \quad \text{and } \sup_{\theta \in \Theta} \left| \frac{1}{n} E \left[\log \left(\frac{\hat{L}_n(\theta)}{\hat{L}_n(\theta_0)} \right) \right] - \ell(\theta; \theta_0) \right| \rightarrow 0$$

for some continuous function $\ell(\theta; \theta_0)$ of θ s.t. for any $\delta > 0$,

$$\sup_{\theta \in \Theta, \|\theta - \theta_0\| \geq \delta} \ell(\theta; \theta_0) < 0, \quad \dots (*)$$

Theorem 3 (consistency of the MLE) Assumption 2 implies

$$\hat{\theta}_n \xrightarrow{p} \theta_0$$

Note

- Assumption 2 can replace Assumption 1 if we only care about consistency of the MLE
- Can verify Assumption 2 for specific cases by checking that the two conditions hold
 - The ULLN holds if Θ is compact, $\log \hat{L}_n(\theta)$ is continuous
 - (*) holds if Θ is compact, $\ell(\theta; \theta_0)$ is continuous on Θ and θ_0 is unique

Assumption 3 As $n \rightarrow \infty$, for each pair $i_1, i_2 \in \{1, \dots, d\}$

$$\underbrace{\sup_{\theta \in \Theta} \left| \frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} - E \left[\frac{\partial^2 \frac{1}{n} \log \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right] \right|}_{\dots(32)} \xrightarrow{p} 0, \quad \text{and } \underbrace{\sup_{\theta \in \Theta} \left| E \left[\frac{\partial^2 \frac{1}{n} \log \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right] - h_{i_1, i_2}(\theta) \right|}_{\dots(33)} \rightarrow 0$$

for some function $h_{i_1, i_2}(\theta)$ that is continuous at θ_0 .

In addition, the (d×d)-matrix $\bar{\mathbf{H}}$ with

$$\underbrace{\bar{\mathbf{H}}_0}_{(d \times d)} = [h_{i_1, i_2}(\theta_0)] \text{ is nonsingular, and } \frac{\partial \frac{1}{\sqrt{n}} \log \hat{L}_n(\theta_0)}{\partial \theta_0^\top} \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{H}}_0)$$

Note: If $\hat{\theta}_n \xrightarrow{p} \theta_0$ (possibly implied by Assumption 2) then Assumption 1, 3 imply that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{H}}_0^{-1})$$

Note: $\bar{\mathbf{H}}_0 = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{H}_n$ where \mathbf{H}_n is the Fisher information

$$\mathbf{H}_n = \text{Var} \left[\frac{\partial \log \hat{L}_n(\theta)}{\partial \theta^\top} \bigg|_{\theta = \theta_0} \right]$$

Note: Verifying (32) implies verification of a ULLN

Note: (33) allows for non-iid data. In the case of iid data we simply have

$$h_{i_1, i_2} = -\frac{1}{n} E \left[\frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta_{i_1} \partial \theta_{i_2}} \right]$$

Theorem 4 (asymptotic Normality of the MLE) Suppose $\hat{\theta}_n \xrightarrow{p} \theta_0$ then Assumption 1, 3 imply $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{H}}_0^{-1})$

Proof By the mean-value theorem for each $i \in \{1, \dots, d\} \exists \hat{\lambda}_i \in (0, 1)$ such that

$$\frac{1}{\sqrt{n}} \frac{\partial \log \hat{L}_n(\theta)}{\partial \theta_i} \bigg|_{\theta = \hat{\theta}_n} = \frac{1}{\sqrt{n}} \frac{\partial \log \hat{L}_n(\theta)}{\partial \theta_i} \bigg|_{\theta = \theta_0} + \frac{1}{n} \frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta \partial \theta_i} \bigg|_{\theta = \theta_0 + \hat{\lambda}_i(\hat{\theta}_n - \theta_0)} \sqrt{n}(\hat{\theta}_n - \theta_0)$$

By FOC for log-likelihood maximization (this requires $\theta_0 \in \text{int} \Theta$ and other conditions of Assumption 1), we have

$$\frac{1}{\sqrt{n}} \frac{\partial \log \hat{L}_n(\theta)}{\partial \theta_i} \bigg|_{\theta = \hat{\theta}_n} = 0, \quad \dots (36)$$

Since Θ is convex (Assumption 1) the mean value satisfies $\theta_0 + \hat{\lambda}_i(\hat{\theta}_n - \theta_0) \in \Theta$, since $\hat{\theta}_n \xrightarrow{p} \theta_0$, the mean value is also consistent for θ_0 and

$$\tilde{\mathbf{H}}_n \equiv \frac{1}{n} \begin{bmatrix} \frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta \partial \theta_1} \bigg|_{\theta = \theta_0 + \hat{\lambda}_1(\hat{\theta}_n - \theta_0)} \\ \vdots \\ \frac{\partial^2 \log \hat{L}_n(\theta)}{\partial \theta \partial \theta_d} \bigg|_{\theta = \theta_0 + \hat{\lambda}_d(\hat{\theta}_n - \theta_0)} \end{bmatrix} \xrightarrow{p} \mathbf{H}_0$$

by the continuous mapping theorem (CMT hereafter) and information identity (implied by Assumption 1).

Since $\bar{\mathbf{H}}_0$ is nonsingular (Assumption 3) then by the CMT

$$\tilde{\mathbf{H}}_n^{-1} \xrightarrow{p} \bar{\mathbf{H}}_0^{-1}, \quad \dots (38)$$

By the original mean value expansion and (36), we have

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = \tilde{\mathbf{H}}_n^{-1} \frac{1}{\sqrt{n}} \frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$$

So by (38) and Slutsky's theorem

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{H}}_0^{-1})$$

(obviously used the assumption that)

$$\frac{1}{\sqrt{n}} \frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \xrightarrow{d} N(\mathbf{0}, \bar{\mathbf{H}}_0)$$

Note: Suppose iid data, i.e. suppose $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ are iid k variate with common density $f(\mathbf{z}; \boldsymbol{\theta}_0)$.

In this case, Assumption 1 implies

$$\frac{1}{n} E \left[\frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = E \left[\frac{\partial \log f(\mathbf{Z}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = \mathbf{0}$$

and

$$\frac{1}{n} \text{Var} \left[\frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = \text{Var} \left[\frac{\partial \log f(\mathbf{Z}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = \bar{\mathbf{H}}_0$$

In this case, it is clear that

$$\frac{1}{\sqrt{n}} \frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial \log f(\mathbf{Z}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{H}}_0)$$

Stylized fact: Subject to regularity conditions (e.g. consistency, Assumption 1 & 3) the MLE is asymptotically efficient over the class of consistent and asymptotically Normal extremum estimators of $\boldsymbol{\theta}_0$.

Consider

$$\begin{aligned} \tilde{\boldsymbol{\theta}}_n &\equiv \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \frac{1}{n} \sum_{j=1}^n g(\mathbf{Z}_j; \boldsymbol{\theta}), \quad \text{i.e. an extremum estimator of} \\ \boldsymbol{\theta}_0 &\equiv \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} E[g(\mathbf{Z}_1, \boldsymbol{\theta})], \quad \text{where } \mathbf{Z}_1, \dots, \mathbf{Z}_n \text{ are iid k vectors} \end{aligned}$$

from an absolutely continuous distribution with density $f(\mathbf{z}; \boldsymbol{\theta}_0)$ where $\boldsymbol{\theta}_0 \in \text{int} \Theta$, where $\Theta \subset \mathbb{R}^d$ and where Θ is convex.

Then $\sqrt{n}(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$ (subject to the right conditions on $g(\cdot, \cdot)$)

$$\text{where } \mathbf{A}_0 = E \left[\frac{\partial^2 g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right]$$

$$\mathbf{B}_0 = E \left[\frac{\partial g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right]$$

Then $\mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1} - \bar{\mathbf{H}}_0^{-1}$ is positive semidefinite, i.e. the MLE standard errors are no larger than the $\tilde{\boldsymbol{\theta}}_n$ standard errors.

To see this: By Assumption 1 and the FOC for

$$\boldsymbol{\theta}_0 = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} E[g(\mathbf{Z}_1, \boldsymbol{\theta})]$$

We have

$$\int \frac{\partial g(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} f(\mathbf{z}; \boldsymbol{\theta}_0) d\mathbf{z} = \mathbf{0}, \quad \dots (49)$$

(obviously) (49) holds regardless of the actual value of $\boldsymbol{\theta}_0$, so

$$\int \frac{\partial g(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} f(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} = \mathbf{0}, \quad \forall \boldsymbol{\theta} \in \Theta, \quad \dots (50)$$

Under the conditions of Assumption 1, we have

$$\begin{aligned} \mathbf{0} &= \int \frac{\partial^2 g(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} f(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} + \int \frac{\partial g(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial f(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} d\mathbf{z} \\ &= \int \frac{\partial^2 g(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} f(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} + \int \frac{\partial g(\mathbf{z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial \log f(\mathbf{z}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} f(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z}, \quad \dots (51) \end{aligned}$$

Let $\boldsymbol{\theta}=\boldsymbol{\theta}_0$, then from (51) we have

$$E \left[\frac{\partial g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial \log f(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = -\mathbf{A}_0, \quad \dots (52)$$

But $\frac{\partial g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$ has mean zero by (49). While $\frac{\partial \log f(\mathbf{Z}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$ also has mean zero by the score identity. Therefore, (52) is equivalent to

$$\operatorname{Cov} \left[\frac{\partial g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}, \frac{\partial \log f(\mathbf{Z}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = -\mathbf{A}_0$$

So by Assumption 3 we have

$$\text{Var} \left[\underbrace{\begin{bmatrix} \frac{\partial g(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \\ \frac{\partial \log f(\mathbf{Z}_1; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \end{bmatrix}}_{2d \times 1} \right] = \underbrace{\begin{bmatrix} \mathbf{B}_0 & -\mathbf{A}_0 \\ -\mathbf{A}_0 & \bar{\mathbf{H}}_0 \end{bmatrix}}_{2d \times 2d}$$

which (by virtue of being a covariance matrix) is positive semidefinite. Therefore

$$\mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1} - \bar{\mathbf{H}}_0^{-1} = [\mathbf{A}_0^{-1} \quad \bar{\mathbf{H}}_0^{-1}] \begin{bmatrix} \mathbf{B}_0 & -\mathbf{A}_0^{-1} \\ -\mathbf{A}_0^{-1} & \bar{\mathbf{H}}_0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_0^{-1} \\ \bar{\mathbf{H}}_0^{-1} \end{bmatrix}$$

is also positive semidefinite.

(This argument also requires one to check that)

$$\frac{1}{\sqrt{n}} \begin{bmatrix} \sum_{j=1}^1 \frac{\partial g(\mathbf{Z}_j, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \\ \frac{\partial \log \hat{\mathbf{L}}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \end{bmatrix} \xrightarrow{d} N_{2d} \left(\mathbf{0}, \begin{bmatrix} \mathbf{B}_0 & -\mathbf{A}_0 \\ -\mathbf{A}_0 & \bar{\mathbf{H}}_0 \end{bmatrix} \right)$$

Estimators based on estimating Equations (i.e. Z-estimators)

- An extremum estimator maximizes a criterion function $\hat{Q}_n(\boldsymbol{\theta})$ globally over a compact set Θ
- This may be impractical
- It may be difficult to maximize $\hat{Q}_n(\boldsymbol{\theta})$ globally, if indeed such a global maximum exists (exception: $\hat{Q}_n(\boldsymbol{\theta})$ known to be globally concave)
- It will be difficult to prove that $\bar{Q}(\boldsymbol{\theta}) \equiv E[\hat{Q}_n(\boldsymbol{\theta})]$ has a unique global maximum on Θ at $\boldsymbol{\theta}_0$
- Specifying Θ as compact may be implausible

e.g. Suppose that

$$E[Y_i | \mathbf{X}_i] = \mathbf{e}^{\mathbf{X}_i^T \boldsymbol{\theta}_0}$$

Why should $\boldsymbol{\theta}_0$ be restricted to lie in a compact subset of \mathbb{R}^d ?

- Proving asymptotic Normality is easiest if $\hat{\boldsymbol{\theta}}_n$ is simply a local maximizer of $\hat{Q}_n(\boldsymbol{\theta})$

For all these reasons, it may be better to consider estimators defined as the solutions to FOCs for maximization of $\hat{Q}_n(\boldsymbol{\theta})$.

$$\text{Let } \hat{Q}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{Z}_i, \boldsymbol{\theta})$$

$$\text{where } \mathbf{Z}_i = [\mathbf{X}_i^T \quad Y_i]^T$$

$$\text{and } \boldsymbol{\theta} \in \mathbb{R}^d$$

Assume $\{\mathbf{Z}_i: i = 1, \dots, n\}$ is iid

$$\text{Let } \bar{Q}(\boldsymbol{\theta}) \equiv E[\hat{Q}_n(\boldsymbol{\theta})] = E[g(\mathbf{Z}_1, \boldsymbol{\theta})]$$

Theorem 1 Assume

- (1) $\Theta \subset \mathbb{R}^d$ is an open set so $\boldsymbol{\theta}_0 \in \text{int}(\Theta)$
- (2) $g(\mathbf{Z}_1, \boldsymbol{\theta})$ is Borel measurable $\forall \boldsymbol{\theta} \in \Theta$, $\frac{\partial g}{\partial \boldsymbol{\theta}^T}$ exists and is continuous $\forall \boldsymbol{\theta} \in N_1(\boldsymbol{\theta}_0)$, where $N_1(\boldsymbol{\theta})$ is an open neighborhood of $\boldsymbol{\theta}_0$ (so $\hat{Q}_n(\boldsymbol{\theta})$ is continuous $\forall \boldsymbol{\theta} \in N_1(\boldsymbol{\theta}_0)$)
- (3) There is an open neighborhood $N_2(\boldsymbol{\theta}_0)$ of $\boldsymbol{\theta}_0$ s.t.

$$\sup_{\boldsymbol{\theta} \in N_2(\boldsymbol{\theta}_0)} |\hat{Q}_n(\boldsymbol{\theta}) - \bar{Q}(\boldsymbol{\theta})| \xrightarrow{P} 0$$

and where $\boldsymbol{\theta}_0$ is the unique (local) maximizer of $\bar{Q}(\boldsymbol{\theta})$ on $N_2(\boldsymbol{\theta}_0)$

Let Θ_n be the set of solutions to $\frac{\partial \hat{Q}_n}{\partial \boldsymbol{\theta}^T} = \mathbf{0}$

(i.e. Θ_n be the set of local extrema of $\hat{Q}_n(\boldsymbol{\theta})$)

(if $\Theta_n = \emptyset$, then define $\Theta_n \equiv \{\emptyset\}$)

Then

$$\forall \varepsilon > 0, \quad \lim_{n \rightarrow \infty} P \left[\inf_{\boldsymbol{\theta} \in \Theta_n} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0) > \varepsilon \right] = \lim_{n \rightarrow \infty} P \left[\inf_{\boldsymbol{\theta} \in \Theta_n} \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|^2 > \varepsilon \right] = 0$$

(Therefore there is a consistent root for the system $\frac{\partial \hat{Q}_n}{\partial \theta^\top} = \mathbf{0}$)

Proof Let $S \subset N_1(\theta_0) \cap N_2(\theta_0)$ be compact. Let $\tilde{\theta}_n$ be the maximizer of $\hat{Q}_n(\theta)$ on S . Then $\tilde{\theta}_n \xrightarrow{p} \theta_0$ (see Theorem 2 in the handout on extremum estimation). But since $\hat{Q}_n(\theta)$ attains a (local) maximum at $\hat{\theta}_n$ with probability approaching 1 as $n \rightarrow \infty$, we have

$$\lim_{n \rightarrow \infty} P[\tilde{\theta}_n \in \Theta_n] = 1$$

A consistent root of $\frac{\partial \hat{Q}_n}{\partial \theta^\top} = \mathbf{0}$ is asymptotically Normal under additional conditions.

Theorem 2 Assume the conditions of Theorem 1 and the following additional conditions

(4) $\frac{\partial^2 \hat{Q}_n}{\partial \theta \partial \theta^\top}$ exists and is continuous on an open convex neighborhood of θ_0

(5)

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial^2 g(\mathbf{Z}_i, \theta)}{\partial \theta \partial \theta^\top} \Big|_{\theta=\tilde{\theta}_n} \xrightarrow{p} \mathbf{A}_0 \equiv E \left[\frac{\partial^2}{\partial \theta \partial \theta^\top} g(\mathbf{Z}_1, \theta) \Big|_{\theta=\theta_0} \right]$$

Where \mathbf{A}_0 is finite and nonsingular for any $\tilde{\theta}_n \xrightarrow{p} \theta_0$

(6)

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial g(\mathbf{Z}_i, \theta)}{\partial \theta^\top} \Big|_{\theta=\theta_0} \xrightarrow{d} N(\mathbf{0}, \mathbf{B}_0)$$

where $\mathbf{B}_0 = E \left[\frac{\partial g(\mathbf{Z}_1, \theta)}{\partial \theta^\top} \Big|_{\theta=\theta_0} \frac{\partial g(\mathbf{Z}_1, \theta)}{\partial \theta} \Big|_{\theta=\theta_0} \right]$

Let $\{\hat{\theta}_n; n = 1, 2, \dots\}$ be a sequence formed by taking one element from Θ_n for each n such that $\hat{\theta}_n \xrightarrow{p} \theta_0$.

Then $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$

Proof

Take a Taylor expansion to get

$$(0 =) \frac{\partial \hat{Q}_n}{\partial \theta^\top} \Big|_{\theta=\hat{\theta}_n} = \frac{\partial \hat{Q}_n}{\partial \theta^\top} \Big|_{\theta=\theta_0} + \frac{\partial^2 \hat{Q}_n}{\partial \theta \partial \theta^\top} \Big|_{\theta=\theta_n^*} (\hat{\theta}_n - \theta_0)$$

Where θ_n^* is a convex combination of $\hat{\theta}_n$ and θ_0 . But $\hat{\theta}_n \in \Theta_n$, so

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = - \left(\frac{\partial^2 \hat{Q}_n}{\partial \theta \partial \theta^\top} \Big|_{\theta=\theta_n^*} \right)^{-1} \sqrt{n} \frac{\partial \hat{Q}_n}{\partial \theta^\top} \Big|_{\theta=\theta_0}$$

But $\boldsymbol{\theta}_n^* \xrightarrow{p} \boldsymbol{\theta}_0$, so by Assumptions (4), (5)

$$\left. \frac{\partial^2 \widehat{Q}_n}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_n^*} \xrightarrow{p} \mathbf{A}_0$$

The conclusion follows from Assumption (6) and Slutsky's theorem.

Note: Assumption (5) in Theorem 2 is usually difficult to verify.

Consider two results using different assumptions that imply Assumption (5).

Let $\mathbf{Z}=[Z_1 \dots Z_n]^\top$. Define $g_n(\boldsymbol{\theta}) \equiv g(\mathbf{Z}, \boldsymbol{\theta})$, where $\boldsymbol{\theta} \in \mathbb{R}^d$, suppose $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$, where Θ open.

Assume that $g(\mathbf{Z}, \boldsymbol{\theta})$ is Borel measurable $\forall \boldsymbol{\theta} \in \Theta$.

Theorem 3 Suppose Θ is convex and that $\frac{\partial g_n}{\partial \boldsymbol{\theta}}$ exists $\forall \boldsymbol{\theta} \in \Theta$. Suppose that $\forall \varepsilon > 0 \exists M_2$ s.t. for each θ_j (where θ_j is the j -th element of $\boldsymbol{\theta}$)

$$P \left[\left| \frac{\partial g_n}{\partial \theta_j} \right| < M_2 \right] \geq 1 - \varepsilon$$

$\forall n \in \{1, 2, \dots\}$ and $\forall \boldsymbol{\theta} \in \Theta$. Then

$$|g_n(\widehat{\boldsymbol{\theta}}_n) - g_n(\boldsymbol{\theta}_0)| \xrightarrow{p} 0, \quad \text{if } \widehat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}_0 \in \Theta$$

Proof

Take a Taylor expansion

$$g_n(\widehat{\boldsymbol{\theta}}_n) = g_n(\boldsymbol{\theta}_0) + \left. \frac{\partial g_n}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_n^*} (\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$$

Where $\boldsymbol{\theta}_n^*$ is between $\widehat{\boldsymbol{\theta}}_n$ and $\boldsymbol{\theta}_0$. By convexity, $\boldsymbol{\theta}_n^* \in \Theta$.

Theorem 4 Suppose \exists an open neighborhood $N(\boldsymbol{\theta}_0)$ of $\boldsymbol{\theta}_0$ such that

$$\sup_{\boldsymbol{\theta} \in N(\boldsymbol{\theta}_0)} |g_n(\boldsymbol{\theta}) - E_{\boldsymbol{\theta}_0}[g_n(\boldsymbol{\theta})]| \xrightarrow{p} 0, \quad (*)$$

Then, $|g_n(\widehat{\boldsymbol{\theta}}_n) - E_{\boldsymbol{\theta}_0}[g_n(\boldsymbol{\theta}_0)]| \xrightarrow{p} 0$ if $\widehat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}_0$ and $E_{\boldsymbol{\theta}_0}[g_n(\boldsymbol{\theta})]$ is continuous at $\boldsymbol{\theta}=\boldsymbol{\theta}_0$.

Proof

Because of the ULLN in (*) and also because $\widehat{\boldsymbol{\theta}}_n \in N(\boldsymbol{\theta}_0)$ with probability approaching 1 as $n \rightarrow \infty$, then $\forall \varepsilon > 0, \forall \delta > 0, \exists M_1$ s.t. $\forall n \geq M_1$,

$$P \left[|g_n(\widehat{\boldsymbol{\theta}}_n) - E_{\boldsymbol{\theta}_0}[g_n(\widehat{\boldsymbol{\theta}}_n)]| \geq \frac{\varepsilon}{2} \right] \leq \frac{\delta}{2}, \quad \dots (†)$$

(exercise: supply details)

Because $E_{\boldsymbol{\theta}_0}[g_n(\boldsymbol{\theta})]$ is continuous at $\boldsymbol{\theta}=\boldsymbol{\theta}_0$, so

$$|E_{\boldsymbol{\theta}_0}[g_n(\widehat{\boldsymbol{\theta}}_n)] - E_{\boldsymbol{\theta}_0}[g_n(\boldsymbol{\theta}_0)]| \xrightarrow{p} 0, \quad \text{by the continuous mapping theorem}$$

So $\forall \varepsilon > 0 \forall \delta > 0, \exists M_2$ s.t. $\forall n \geq M_2$,

$$P \left[\left| E_{\theta_0} [g_n(\hat{\theta}_n)] - E_{\theta_0} [g_n(\theta_0)] \right| \geq \frac{\varepsilon}{2} \right] \leq \frac{\delta}{2}, \quad \dots (††)$$

From (†), (††), we have $\forall n \geq \max\{M_1, M_2\}$

$$P[|g_n(\hat{\theta}_n) - g_n(\theta_0)| \leq \varepsilon] \geq 1 - \delta$$

Note: $E_{\theta_0} [g_n(\hat{\theta}_n)] \approx \frac{1}{n} \sum_{i=1}^n E_{\theta_0} [g_n(\hat{\theta}_n) | \hat{\theta}_n]$

Exercise: Use Theorem 1, 2 in this lecture to rederive the consistency and asymptotic Normality of estimators defined as solutions of

$$\frac{\partial}{\partial \theta^\top} \log \hat{L}_n(\theta) = \mathbf{0}$$

Where $\hat{L}_n(\theta)$ satisfies L1–L3.

Quasi-Maximum Likelihood (also known as Pseudo-Maximum Likelihood)

- In practice, one can never be sure that an MLE actually maximizes a random function satisfying L1–L3
- This can be serious because, misspecification of $\hat{L}_n(\theta)$ may cause the score identity to fail, leading the MLE to be inconsistent
- On the other hand, a quasi-MLE, defined as the maximizer of a possibly misspecified log-likelihood will generally be consistent for the pseudo-true value, i.e. for

$$\bar{\theta}_0 \equiv \operatorname{argmax}_{\theta \in \Theta} E_{\theta_0} \left[\frac{1}{n} \log \hat{L}_n(\theta) \right], \quad \dots (*)$$

where the expectation is taken with respect to the time data density indexed by θ_0 . From (*), we see that if $\hat{L}_n(\theta)$ is correctly satisfied, then $\bar{\theta}_0 = \theta_0$; otherwise the quasi-MLE is inconsistent.

Huber (1962), White (1982) showed that a quasi-MLE given as

$$\hat{\boldsymbol{\theta}}_{\text{QML}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta})$$

is consistent for the pseudo-true value, i.e. for

$$\bar{\boldsymbol{\theta}}_0 = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} E_{\boldsymbol{\theta}_0} \left[\frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta}) \right]$$

(where the expectation is taken w.r.t. the true joint data density indexed by $\boldsymbol{\theta}_0$)

In fact, we have

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{\text{QML}} - \bar{\boldsymbol{\theta}}_0) \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{A}}_0^{-1} \bar{\mathbf{B}}_0 \bar{\mathbf{A}}_0^{-1})$$

where

$$\bar{\mathbf{A}}_0 = E_{\boldsymbol{\theta}_0} \left[\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log \hat{L}_n(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right], \quad \bar{\mathbf{B}}_0 = E_{\boldsymbol{\theta}_0} \left[\frac{\partial}{\partial \boldsymbol{\theta}^\top} \log \hat{L}_n(\boldsymbol{\theta}) \frac{\partial}{\partial \boldsymbol{\theta}} \log \hat{L}_n(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right]$$

If $\bar{\boldsymbol{\theta}}_0 = \boldsymbol{\theta}_0$, i.e. the likelihood is correctly specified, then by the information identity

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{\text{QML}} - \boldsymbol{\theta}_0) \xrightarrow{d} N_d(\mathbf{0}, \bar{\mathbf{A}}_0^{-1})$$

How to interpret $\bar{\boldsymbol{\theta}}_0$?

White (1982) proposed the interpretation.

- Let $f(\mathbf{z}; \boldsymbol{\theta})$ denote the assumed joint density of $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ (where $\boldsymbol{\theta} \in \Theta$)
- Let $g(\mathbf{z})$ denote the true but unknown density of $\mathbf{Z}_1, \dots, \mathbf{Z}_n$

Then the Kullback–Leibler information criterion (KLIC) as

$$\text{KL}(g; f) \equiv E_g \left[\log \left(\frac{g(\mathbf{z})}{f(\mathbf{z}; \boldsymbol{\theta})} \right) \right] = E_g[\log g(\mathbf{z})] - E_g[\log f(\mathbf{z}; \boldsymbol{\theta})]$$

(the expectation w.r.t. $g(\mathbf{z})$)

The KLIC is minimized at 0 if $\exists \boldsymbol{\theta}_0$ s.t.

$$g(\mathbf{z}) = f(\mathbf{z}; \boldsymbol{\theta}_0), \quad \forall \mathbf{z}$$

i.e. the model $\{f(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\theta} \in \Theta\}$ is correctly specified White (1982) showed that the QMLE $\hat{\boldsymbol{\theta}}_{\text{QML}}$ minimizes

$$\text{KL}(g; f(\cdot; \boldsymbol{\theta})), \quad \text{over } \boldsymbol{\theta} \in \Theta$$

because

$$\frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta}) \xrightarrow{p} E_g[\log f(\mathbf{z}; \boldsymbol{\theta})], \quad \text{under suitable conditions}$$

and $E_g[\log f(\mathbf{z}; \boldsymbol{\theta})]$ is maximized over Θ by $\bar{\boldsymbol{\theta}}_0$, i.e. $\bar{\boldsymbol{\theta}}_0$ minimizes $\text{KL}(g; f(\cdot; \boldsymbol{\theta}))$ over $\boldsymbol{\theta} \in \Theta$.

Therefore, $\hat{\boldsymbol{\theta}}_{\text{QML}}$ is consistent for the minimizer of $\text{KL}(g; f(\cdot; \boldsymbol{\theta}))$, which “may be useful” if $\text{KL}(g; f(\cdot; \boldsymbol{\theta}))$ is small.

Next: Nonlinear least squares

Suppose we observe $[\mathbf{X}_i^\top Y_i]^\top$ ($i=1, \dots, n$) iid, where \mathbf{X}_i is k -variate.

Assume $\exists \boldsymbol{\beta} \in \mathbb{R}^d$ and a function $g: \mathbb{R}^k \times \mathbb{R}^d \rightarrow \mathbb{R}$ s.t. $E[Y_i | \mathbf{X}_i] = g(\mathbf{X}_i; \boldsymbol{\beta}_0)$ w.p. 1, i.e. a nonlinear regression model.

In this case, a natural estimator of $\boldsymbol{\beta}_0$ is the nonlinear least squares (NLLS) estimator.

$$\tilde{\boldsymbol{\beta}}_{\text{NLLS}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^d}{\operatorname{argmax}} \hat{Q}_n(\boldsymbol{\beta})$$

where

$$\hat{Q}_n(\boldsymbol{\beta}) = -\frac{1}{2n} \sum_{i=1}^n (Y_i - g(\mathbf{X}_i, \boldsymbol{\beta}))^2$$

Somewhat confusingly, “nonlinear least squares” usually refers to the following Z-estimator of $\boldsymbol{\beta}_0$:

$$\left. \frac{\partial \hat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \right|_{\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}_{\text{NLLS}}} = \mathbf{0}$$

Under the conditions for Z-estimation $\hat{\boldsymbol{\beta}}_{\text{NLLS}}$ is \sqrt{n} -consistent with

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{NLLS}} - \boldsymbol{\beta}_0) \xrightarrow{d} N_d(\mathbf{0}, \mathbf{A}_0^{-1} \mathbf{B}_0 \mathbf{A}_0^{-1})$$

where

$$\mathbf{A}_0 = -E \left[\left. \frac{\partial g(\mathbf{X}_1, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \frac{\partial g(\mathbf{X}_1, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta} = \boldsymbol{\beta}_0} \right], \quad \mathbf{B}_0 = E \left[u_1^2 \left. \frac{\partial g(\mathbf{X}_1, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \frac{\partial g(\mathbf{X}_1, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta} = \boldsymbol{\beta}_0} \right]$$

where $u_1 = Y_1 - g(\mathbf{X}_1, \boldsymbol{\beta}_0)$

(Exercise: Fill in the details and give specific assumptions)

e.g. Suppose linear regression, i.e. $k=d$,

$$g(\mathbf{x}, \boldsymbol{\beta}_0) = \mathbf{x}^\top \boldsymbol{\beta}_0, \quad \forall \mathbf{x} \in \mathbb{R}^k$$

Then $\hat{\boldsymbol{\beta}}_{\text{NLLS}} = \hat{\boldsymbol{\beta}}_{\text{OLS}}$ and $\mathbf{A}_0 = -E[\mathbf{X}_1 \mathbf{X}_1^\top]$, $\mathbf{B}_0 = E[u_1^2 \mathbf{X}_1 \mathbf{X}_1^\top]$.

e.g. Suppose $E[Y_i | \mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}$ ($k=d$). Then,

$$\left. \frac{\partial \hat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \right|_{\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}_{\text{NLLS}}} = \mathbf{0} \Leftrightarrow \frac{1}{n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_{\text{NLLS}}}) e^{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_{\text{NLLS}}} \mathbf{X}_i = \mathbf{0}$$

i.e. no closed-form expression for $\hat{\boldsymbol{\beta}}_{\text{NLLS}}$

Also,

$$\mathbf{A}_0 = -E[e^{2\mathbf{X}_1^\top \boldsymbol{\beta}_0} \mathbf{X}_1 \mathbf{X}_1^\top], \quad \mathbf{B}_0 = E[u_1^2 e^{2\mathbf{X}_1^\top \boldsymbol{\beta}_0} \mathbf{X}_1 \mathbf{X}_1^\top]$$

In this case, the asymptotic variance of $\hat{\beta}_{NLLS}$ is consistently estimated as

$$\hat{V}_n[\hat{\beta}_{NLLS}] = \frac{1}{n} \hat{A}_n^{-1} \hat{B}_n \hat{A}_n^{-1}$$

where

$$\hat{A}_n = -\frac{1}{n} \sum_{i=1}^n e^{2\mathbf{X}_i^T \hat{\beta}_{NLLS}} \mathbf{X}_i \mathbf{X}_i^T, \quad \hat{B}_n = \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 e^{2\mathbf{X}_i^T \hat{\beta}_{NLLS}} \mathbf{X}_i \mathbf{X}_i^T, \quad \hat{u}_i = Y_i - e^{\mathbf{X}_i^T \hat{\beta}_{NLLS}}$$

Note: $\hat{V}_n[\hat{\beta}_{NLLS}]$ is heteroskedasticity-consistent, i.e. $E[u_i^2 | \mathbf{X}_i]$ can vary with i .

(On the other hand, if $E[u_i^2 | \mathbf{X}_i] = \sigma_u^2$ w.p. 1 (i.e. homoskedastic), then can estimate the NLLS asymptotic variance via $s_n^2 \hat{A}_n^{-1}$ where)

$$s_n^2 = \frac{1}{n-k} \sum_{i=1}^n \hat{u}_i^2$$

In general, let

$$\hat{A}_n = -\frac{1}{n} \sum_{i=1}^n \frac{\partial g(\mathbf{X}_i, \beta)}{\partial \beta^T} \frac{\partial g(\mathbf{X}_i, \beta)}{\partial \beta} \Big|_{\beta=\hat{\beta}_{NLLS}}$$

$$\hat{B}_n = \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \frac{\partial g(\mathbf{X}_i, \beta)}{\partial \beta^T} \frac{\partial g(\mathbf{X}_i, \beta)}{\partial \beta} \Big|_{\beta=\hat{\beta}_{NLLS}}$$

where $\hat{u}_i = Y_i - g(\mathbf{X}_i, \hat{\beta}_{NLLS})$

Convenient to have an expression for the heteroskedasticity-consistent asymptotic variance estimator $\hat{V}_n[\hat{\beta}_{NLLS}]$ in terms of matrices

$$\text{Let } \hat{\Omega}_n \equiv \text{diag}[\hat{u}_1^2, \dots, \hat{u}_n^2]$$

$$\frac{\partial \mathbf{g}}{\partial \beta} = \left[\frac{\partial g(\mathbf{X}_i, \beta)}{\partial \beta_j} \right]_{i=1, j=1}^{n, d}, \quad \hat{G}_n \equiv \frac{\partial \mathbf{g}}{\partial \beta} \Big|_{\beta=\hat{\beta}_{NLLS}}$$

Then,

$$\hat{A}_n = -\frac{1}{n} \hat{G}_n^T \hat{G}_n, \quad \hat{B}_n = \frac{1}{n} \hat{G}_n^T \hat{\Omega}_n \hat{G}_n$$

So

$$\hat{V}_n[\hat{\beta}_{NLLS}] = \frac{1}{n} \hat{A}_n^{-1} \hat{B}_n \hat{A}_n^{-1} = (\hat{G}_n^T \hat{G}_n)^{-1} \hat{G}_n^T \hat{\Omega}_n \hat{G}_n (\hat{G}_n^T \hat{G}_n)^{-1}$$

Note: One often uses $\hat{B}_n = \frac{1}{n-d} \hat{G}_n^T \hat{\Omega}_n \hat{G}_n$, which makes

$$\hat{V}_n[\hat{\beta}_{NLLS}] = \frac{n}{n-d} (\hat{G}_n^T \hat{G}_n)^{-1} \hat{G}_n^T \hat{\Omega}_n \hat{G}_n (\hat{G}_n^T \hat{G}_n)^{-1}$$

Consider feasible generalized nonlinear least squares (FGNLLS)

Suppose we model heteroskedasticity as

$$E[\mathbf{u}\mathbf{u}^\top|\mathbf{X}] = \mathbf{\Omega}(\boldsymbol{\gamma}_0), \quad \mathbf{u} = [u_1 \quad \cdots \quad u_n]^\top, \quad \mathbf{X} = [\mathbf{X}_1 \quad \cdots \quad \mathbf{X}_n]^\top$$

Suppose $\exists \hat{\boldsymbol{\gamma}}_n \xrightarrow{p} \boldsymbol{\gamma}_0$. The FGNLLS estimator of $\boldsymbol{\beta}$ maximizes the weighted sum of squares

$$\hat{Q}_n(\boldsymbol{\beta}) = -\frac{1}{2n}(\mathbf{y} - \mathbf{g}(\boldsymbol{\beta}))^\top [\mathbf{\Omega}(\hat{\boldsymbol{\gamma}}_n)]^{-1}(\mathbf{y} - \mathbf{g}(\boldsymbol{\beta}))$$

where $\mathbf{y} = [y_1 \quad \cdots \quad y_n]^\top$, $\mathbf{g}(\boldsymbol{\beta}) = [g(\mathbf{X}_1, \boldsymbol{\beta}) \quad \cdots \quad g(\mathbf{X}_n, \boldsymbol{\beta})]^\top$

(Exercise: Specify appropriate assumptions such that)

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{FGNLLS}} - \boldsymbol{\beta}_0) \xrightarrow{d} N_d\left(\mathbf{0}, \left(\text{plim} \frac{1}{n} \hat{\mathbf{G}}_n^\top \mathbf{\Omega}(\hat{\boldsymbol{\gamma}}) \hat{\mathbf{G}}_n\right)^{-1}\right)$$

where

$$\hat{\mathbf{G}}_n \equiv \left. \frac{\partial}{\partial \boldsymbol{\beta}} \mathbf{g}(\boldsymbol{\beta}) \right|_{\boldsymbol{\beta} = \hat{\boldsymbol{\beta}}_{\text{FGNLLS}}}$$

Note: FGNLLS delivers smaller standard errors than available using the heteroskedasticity-consistent NLLS asymptotic variance estimator, but at the price of requiring a model for $E[\mathbf{u}\mathbf{u}^\top|\mathbf{X}]$.

FGNLLS asymptotic variance estimate

$$\left(\hat{\mathbf{G}}_n^\top (\mathbf{\Omega}(\hat{\boldsymbol{\gamma}}_n))^{-1} \hat{\mathbf{G}}_n\right)^{-1}$$

is inconsistent.

Weighted Nonlinear Least Squares (WNLLS) can avoid having to specify a model for $E[\mathbf{u}\mathbf{u}^\top|\mathbf{X}]$ by adopting a working error variance matrix. In particular, let $\boldsymbol{\Sigma}_0 \equiv \boldsymbol{\Sigma}(\boldsymbol{\gamma}_0)$ be a possible model for $E[\mathbf{u}\mathbf{u}^\top|\mathbf{X}]$ (i.e. $\boldsymbol{\Sigma}$ is a working error variance matrix). Assume

$$\hat{\boldsymbol{\gamma}}_n \xrightarrow{p} \boldsymbol{\gamma}_0, \quad \text{Let } \hat{\boldsymbol{\Sigma}}_n = \boldsymbol{\Sigma}(\hat{\boldsymbol{\gamma}}_n)$$

The weighted nonlinear least squares (WNLLS) estimator of $\boldsymbol{\beta}_0$ is

$$\hat{\boldsymbol{\beta}}_{\text{WNLLS}} \equiv \underset{\boldsymbol{\beta} \in \mathbb{R}^d}{\text{argmax}} \hat{Q}_n(\boldsymbol{\beta}), \quad \hat{Q}_n(\boldsymbol{\beta}) = -\frac{1}{2n}(\mathbf{y} - \mathbf{g}(\boldsymbol{\beta}))^\top \hat{\boldsymbol{\Sigma}}_n^{-1}(\mathbf{y} - \mathbf{g}(\boldsymbol{\beta}))$$

Exercise: Specify assumptions such that $\hat{\boldsymbol{\beta}}_{\text{WNLLS}}$ is \sqrt{n} -consistent with

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{\text{WNLLS}} - \boldsymbol{\beta}_0) \xrightarrow{d} N_d\left(\mathbf{0}, \text{plim } n(\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1} \hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\boldsymbol{\Omega}}_n \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n (\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1}\right)$$

Where $\hat{\boldsymbol{\Omega}}_n = \text{diag}[\hat{u}_1^2 \quad \cdots \quad \hat{u}_n^2]$

Note: If $\boldsymbol{\Sigma}(\boldsymbol{\gamma}_0) = \mathbf{\Omega} \equiv E[\mathbf{u}\mathbf{u}^\top|\mathbf{X}]$, then the WNLLS asymptotic variance will be close to the FGNLLS asymptotic variance when n is large.

Alarm: Midterm November 3, 2016

NLLS with a working error variance matrix

Let $\Sigma_0 \equiv \Sigma(\gamma_0)$ be a possible model for $E[\mathbf{u}\mathbf{u}^\top|\mathbf{X}]$. Assume $\hat{\gamma}_n \xrightarrow{p} \gamma_0$. Let $\hat{\Sigma}_n \equiv \Sigma(\hat{\gamma}_n)$.
Weighted NLLS is given by

$$\hat{\beta}_{\text{WNLLS}} \equiv \underset{\beta \in \mathbb{R}^d}{\operatorname{argmax}} \hat{Q}_n(\beta)$$

where

$$\hat{Q}_n(\beta) = -\frac{1}{2n} (\mathbf{y} - \mathbf{g}(\beta))^\top \hat{\Sigma}_n^{-1} (\mathbf{y} - \mathbf{g}(\beta))$$

Then under appropriate assumptions,

$$\sqrt{n}(\hat{\beta}_{\text{WNLLS}} - \beta_0) \xrightarrow{d} N_d \left(0, \operatorname{plim}_n (\hat{\mathbf{G}}_n^\top \hat{\Sigma}_n^{-1} \hat{\mathbf{G}}_n)^{-1} \hat{\mathbf{G}}_n^\top \hat{\Sigma}_n^{-1} \hat{\Omega}_n \hat{\Sigma}_n^{-1} \hat{\mathbf{G}}_n (\hat{\mathbf{G}}_n^\top \hat{\Sigma}_n^{-1} \hat{\mathbf{G}}_n)^{-1} \right)$$

where $\hat{\Omega}_n = \operatorname{diag}[\hat{u}_1^2, \dots, \hat{u}_n^2]$

e.g. Consider WNLLS with a working model of multiplicative heteroskedasticity, i.e.

$$\sigma_i^2 \equiv E[u_i^2 | \mathbf{x}_i] = e^{\mathbf{z}_i^\top \gamma_0}, \quad \text{for each } i \in \{1, \dots, n\}$$

where \mathbf{z}_i is a subvector of \mathbf{x}_i

Then

$$\Sigma_0 = \operatorname{diag} [e^{\mathbf{z}_i^\top \gamma_0} : i = 1, \dots, n]$$

and $\hat{\Sigma}_n = \operatorname{diag} [e^{\mathbf{z}_i^\top \hat{\gamma}_n} : i = 1, \dots, n]$

where $\hat{\gamma}_n$ can be obtained by NLLS of the squared NLLS residuals $(y_i - g(\mathbf{x}_i, \hat{\beta}_{\text{NLLS}}))^2$ on $e^{\mathbf{z}_i^\top \gamma}$.
We have

$$\Sigma_0^{-1} = \operatorname{diag} [e^{-\mathbf{z}_i^\top \gamma_0} : i = 1, \dots, n]$$

And the WNLLS estimator maximizes

$$\hat{Q}_n(\beta) = \frac{1}{2n} \sum_{i=1}^n \frac{(y_i - g(\mathbf{x}_i, \beta))^2}{e^{\mathbf{z}_i^\top \hat{\gamma}_n}}$$

The asymptotic variance matrix of $\hat{\boldsymbol{\beta}}_{\text{WNLLS}}$ is consistently estimated by

$$\begin{aligned} (\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1} \hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\boldsymbol{\Omega}}_n \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n (\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1} &= \left(\sum_{i=1}^n \frac{1}{\hat{\sigma}_i^2} \frac{\partial}{\partial \boldsymbol{\beta}^\top} g(\mathbf{x}_i, \boldsymbol{\beta}) \frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}_i, \boldsymbol{\beta}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{WNLLS}}} \right)^{-1} \\ &\times \left(\sum_{i=1}^n \frac{\hat{u}_i^2}{\hat{\sigma}_i^4} \frac{\partial}{\partial \boldsymbol{\beta}^\top} g(\mathbf{x}_i, \boldsymbol{\beta}) \frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}_i, \boldsymbol{\beta}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{WNLLS}}} \right) \\ &\times \left(\sum_{i=1}^n \frac{1}{\hat{\sigma}_i^2} \frac{\partial}{\partial \boldsymbol{\beta}^\top} g(\mathbf{x}_i, \boldsymbol{\beta}) \frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}_i, \boldsymbol{\beta}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{WNLLS}}} \right)^{-1} \end{aligned}$$

where

$$\hat{\sigma}_i^2 = \mathbf{e}^{\mathbf{z}_i^\top} \hat{\mathbf{y}}_n, \quad \hat{u}_i = y_i - g(\mathbf{x}_i, \hat{\boldsymbol{\beta}}_{\text{WNLLS}})$$

Note: If one assumes that

$$\boldsymbol{\Sigma}_0 = \boldsymbol{\Omega}_0 \equiv \mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}]$$

then WNLLS coincides with FGNLLS and the asymptotic variance matrix is consistently estimated by

$$(\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1} = \left(\sum_{i=1}^n \frac{1}{\hat{\sigma}_i^2} \frac{\partial}{\partial \boldsymbol{\beta}^\top} g(\mathbf{x}_i, \boldsymbol{\beta}) \frac{\partial}{\partial \boldsymbol{\beta}} g(\mathbf{x}_i, \boldsymbol{\beta}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{\text{WNLLS}}} \right)^{-1}$$

e.g. Suppose $y|\mathbf{x} \sim \text{exponential}(e^{\mathbf{x}^\top \boldsymbol{\beta}_0})$, i.e. y given \mathbf{x} has a continuous distribution with density

$$f_{Y|\mathbf{x}}(y; \boldsymbol{\beta}_0) = \begin{cases} e^{\mathbf{x}^\top \boldsymbol{\beta}_0} e^{-e^{\mathbf{x}^\top \boldsymbol{\beta}_0} y}, & y > 0 \\ 0, & y \leq 0 \end{cases}$$

Then $\mathbb{E}[Y|\mathbf{x}] = e^{-\mathbf{x}^\top \boldsymbol{\beta}_0}$. Suppose we observe $[\mathbf{x}_i^\top y_i]^\top$ ($i=1, \dots, n$) iid. In this case, the OLS estimator of $\boldsymbol{\beta}_0$ is inconsistent (exercise: verify this).

The MLE is obtained as follows.

$$\log f_{Y|\mathbf{x}}(y; \boldsymbol{\beta}_0) = \mathbf{x}^\top \boldsymbol{\beta}_0 - e^{\mathbf{x}^\top \boldsymbol{\beta}_0} y, \quad \text{if } y > 0$$

So

$$\begin{aligned} \hat{L}_n(\boldsymbol{\beta}) &\equiv \frac{1}{n} \sum_{i=1}^n \log f_{Y|\mathbf{x}}(y_i; \boldsymbol{\beta}) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^\top \boldsymbol{\beta} - \frac{1}{n} \sum_{i=1}^n e^{\mathbf{x}_i^\top \boldsymbol{\beta}} y_i \\ \Rightarrow \frac{\partial \hat{L}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n e^{\mathbf{x}_i^\top \boldsymbol{\beta}} y_i \mathbf{x}_i \end{aligned}$$

The FOCs for likelihood maximization are

$$\frac{\partial \hat{L}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} = \mathbf{0} \Leftrightarrow \frac{1}{n} \sum_{i=1}^n (1 - e^{\mathbf{x}_i^\top \boldsymbol{\beta}} y_i) \mathbf{x}_i = \mathbf{0}$$

which has no closed-form solution (exercise: check the second-order condition, give a sufficient condition for $\frac{\partial^2 \hat{L}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top}$ to be negative definite).

Since $E[Y|\mathbf{x}] = e^{-\mathbf{x}^\top \boldsymbol{\beta}_0}$, we can write $Y = e^{-\mathbf{x}^\top \boldsymbol{\beta}_0} + u$, where u is a random error with $E[u|\mathbf{x}] = 0$. Note that

$$E[u^2|\mathbf{x}] = \text{Var}[u|\mathbf{x}] = \text{Var}[y|\mathbf{x}] = e^{-2\mathbf{x}^\top \boldsymbol{\beta}_0}$$

So (conditionally) heteroskedastic

$\boldsymbol{\beta}_0$ can be consistently estimated by NLLS by maximizing

$$\hat{Q}_{\text{NLLS}}(\boldsymbol{\beta}) = -\frac{1}{2n} \sum_{i=1}^n (y_i - e^{-\mathbf{x}_i^\top \boldsymbol{\beta}})^2$$

So

$$\frac{\partial \hat{Q}_{\text{NLLS}}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} = -\frac{1}{n} \sum_{i=1}^n (y_i - e^{-\mathbf{x}_i^\top \boldsymbol{\beta}}) e^{-\mathbf{x}_i^\top \boldsymbol{\beta}} \mathbf{x}_i \neq \frac{\partial \hat{L}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}$$

Therefore $\hat{\boldsymbol{\beta}}_{\text{NLLS}} \neq \hat{\boldsymbol{\beta}}_{\text{MLE}}$.

Now consider WNLLS. Choose an error variance proportional to the mean. In particular, set

$$\boldsymbol{\Sigma}_0 = \text{diag}[e^{-\mathbf{x}_i^\top \boldsymbol{\beta}_0} : i = 1, \dots, n]$$

and $\hat{\boldsymbol{\beta}}_{\text{WNLLS}}$ to maximize

$$\begin{aligned} \hat{Q}_{\text{WNLLS}}(\boldsymbol{\beta}) &= -\frac{1}{2n} \sum_{i=1}^n \frac{(y_i - e^{-\mathbf{x}_i^\top \boldsymbol{\beta}})^2}{e^{-\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{\text{NLLS}}}} \\ \Rightarrow \frac{\partial \hat{Q}_{\text{WNLLS}}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} &= -\frac{1}{n} \sum_{i=1}^n (y_i - e^{-\mathbf{x}_i^\top \boldsymbol{\beta}}) e^{-\mathbf{x}_i^\top \boldsymbol{\beta}} e^{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{\text{NLLS}}} \mathbf{x}_i \\ \Rightarrow \frac{\partial \hat{Q}_{\text{WNLLS}}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} &\approx -\frac{1}{n} \sum_{i=1}^n (y_i e^{\mathbf{x}_i^\top \boldsymbol{\beta}_0} - 1) \mathbf{x}_i = -\frac{\partial \hat{L}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \Big|_{\boldsymbol{\beta}=\boldsymbol{\beta}_0} \end{aligned}$$

when n is large (given the consistency of $\hat{\boldsymbol{\beta}}_{\text{NLLS}}$).

Therefore, $\hat{\boldsymbol{\beta}}_{\text{FGNLLS}} \approx \hat{\boldsymbol{\beta}}_{\text{MLE}}$ when n is large, i.e. $\hat{\boldsymbol{\beta}}_{\text{FGNLLS}}$ is asymptotically efficient assuming a correctly specified likelihood.

Consider some simulations.

Suppose $Y|x \sim \text{exponential}(e^{\beta_1 + \beta_2 x})$, where $X \sim N(1, 1^2)$ and $\beta_1=2, \beta_2=-1$. Draw $[x_i \ y_i]^\top$ ($i=1, \dots, 10,000$) from this model. Find $\bar{y}=0.62$ with standard deviation $s_y=1.29$.

1. OLS

$$\hat{\beta}_{1,OLS} = -0.0093, \quad \hat{\beta}_{2,OLS} = 0.6198, \quad \text{inconsistent}$$

2. MLE

$$\hat{\beta}_{1,MLE} = \frac{1.9829}{(0.0144)}, \quad \hat{\beta}_{2,MLE} = \frac{-0.9896}{(0.0099)}$$

The MLE standard errors are derived from

$$\frac{1}{n} \hat{\mathbf{H}}_n^{-1}(\hat{\boldsymbol{\beta}}_{MLE}), \quad \text{where } n = 10,000 \text{ and}$$

$$\hat{\mathbf{H}}_n(\hat{\boldsymbol{\beta}}_{MLE}) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \log f_{Y|x}(y_i; \boldsymbol{\beta}) \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{MLE}}$$

3. NLLS

$$\hat{\beta}_{1,NLLS} = \frac{1.8876}{(0.1421)}, \quad \hat{\beta}_{2,NLLS} = \frac{0.9575}{(0.0612)}$$

The NLLS standard errors come from the heteroskedasticity-consistent estimate

$$\hat{\mathbf{V}}[\hat{\boldsymbol{\beta}}_{NLLS}] = \frac{1}{n} \hat{\mathbf{A}}_n^{-1} \hat{\mathbf{B}}_n \hat{\mathbf{A}}_n^{-1}$$

$$= (\hat{\mathbf{G}}_n^\top \hat{\mathbf{G}}_n)^{-1} \hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Omega}}_n \hat{\mathbf{G}}_n (\hat{\mathbf{G}}_n^\top \hat{\mathbf{G}}_n)^{-1}$$

where

$$\hat{\mathbf{A}}_n = -\frac{1}{n} \sum_{i=1}^n \frac{\partial g(x_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \frac{\partial g(x_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{NLLS}}, \quad \hat{\mathbf{B}}_n = \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2 \frac{\partial g(x_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \frac{\partial g(x_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_{NLLS}}$$

$$\hat{u}_i = y_i - g(x_i, \hat{\boldsymbol{\beta}}_{NLLS}), \quad g(x_i, \boldsymbol{\beta}) = e^{-\beta_1 - \beta_2 x_i}$$

4. WNLLS with the working error variance

$$\hat{\sigma}_i^2 = e^{-\hat{\beta}_{1,NLLS} - \hat{\beta}_{2,NLLS} x_i}$$

we get

$$\hat{\beta}_{1,WNLLS} = \frac{1.9906}{(0.0359)}, \quad \hat{\beta}_{2,WNLLS} = \frac{-0.9961}{(0.0224)}$$

The WNLLS standard errors come from the matrix

$$(\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1} \hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\boldsymbol{\Omega}}_n \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n (\hat{\mathbf{G}}_n^\top \hat{\boldsymbol{\Sigma}}_n^{-1} \hat{\mathbf{G}}_n)^{-1}$$

where

$$\hat{\boldsymbol{\Sigma}}_n = \text{diag}[\hat{\sigma}_i^2: i = 1, \dots, n], \quad \hat{\boldsymbol{\Omega}}_n = \text{diag}[\hat{u}_i^2: i = 1, \dots, n]$$

5. FGNLLS

$$\hat{\beta}_{1,FGNLLS} = \frac{1.9840}{(0.0146)}, \quad \hat{\beta}_{2,FGNLLS} = \frac{-0.9907}{(0.0101)}$$

FGNLLS in this case is basically equal to MLE in terms of performance.

The FGNLLS standard errors come from

$$(\hat{G}_n^T \hat{\Sigma}_n^{-1} \hat{G}_n)^{-1}$$

SAS Program for the Simulation

```

/*****
Replication of Table 5.7 from Cameron and Trivedi "Microeconometrics" (p. 161)
Written by Junyong Kim, University of Wisconsin-Milwaukee on October 25, 2016
*****/
resetline;
ods html close;
ods graphics off;
ods listing;

/*****
1. Generate 10,000 observations from X~Normal(1,1), Y|x~Exponential(exp(2-1*x))
*****/
data _01;
  do I=1 to 10000;
    X=1+rannor(1);
    Y=-log(1-ranuni(2))/exp(2-X);
    output;
  end;
run;

/*****
2. Compute sample mean and standard deviation of yi
*****/
proc means mean std maxdec=2;
  var Y;
run;

/*****
3. Compute OLS, ML, NLS, WNLS, FGNLS estimates from (yi,xi)
*****/
proc iml;
  use _01;
  read all var{Y} into Y;
  read all var{X} into X;
  X2=j(nrow(X),1,1)||X;
  /*****
  3.1. Compute OLS estimates (from the closed-form solution)
  *****/
  BETA_OLS=inv(X2`*X2)*X2`*Y;
  U_OLS=Y-X2*BETA_OLS;
  S2_OLS=U_OLS`*U_OLS/(nrow(X2)-ncol(X2));
  VAR_OLS=S2_OLS*inv(X2`*X2);
  STE_OLS=sqrt(vecdiag(VAR_OLS));
  RSQ_OLS=1-U_OLS`*U_OLS/((Y-mean(Y))`*(Y-mean(Y)));
  /*****
  3.2. Compute ML estimates (by maximizing the log-likelihood with Quasi-Newton method)
  *****/
  start LNHAT(BETA) global(X,Y);
    LNHAT=0;
    do I=1 to nrow(Y);
      LNHAT=LNHAT+BETA[1]+BETA[2]*X[I]-exp(BETA[1]+BETA[2]*X[I])*Y[I];
    end;
    return(LNHAT);
  finish LNHAT;
  call nlpqn(RC1,BETA_MLE,"LNHAT",BETA_OLS,{1,0});

```

```

BETA_MLE=BETA_MLE`;
H_MLE={0 0,0 0};
do I=1 to nrow(Y);
    H_MLE=H_MLE+exp(BETA_MLE[1]+BETA_MLE[2]*X[I])*Y[I]*(1/X[I])*(1|X[I])/nrow(Y);
end;
VAR_MLE=inv(H_MLE)/nrow(Y);
STE_MLE=sqrt(vecdiag(VAR_MLE));
/*****
3.3. Compute Nonlinear Least Squares estimates (by maximizing Q with QN method)
*****/
start QNLLS(BETA) global(X,Y);
    QNLLS=0;
    do I=1 to nrow(Y);
        QNLLS=QNLLS-(Y[I]-exp(-BETA[1]-BETA[2]*X[I]))**2/(2*nrow(Y));
    end;
    return(QNLLS);
finish QNLLS;
call nlpqn(RC2,BETA_NLLS,"QNLLS",BETA_OLS,{1,0});
BETA_NLLS=BETA_NLLS`;
A_NLLS=0;
B_NLLS=0;
do I=1 to nrow(Y);
    A_NLLS=A_NLLS-exp(-2*BETA_NLLS[1]-2*BETA_NLLS[2]*X[I])*(1/X[I])*(1|X[I])/nrow(Y);
    B_NLLS=B_NLLS+(Y[I]-exp(-BETA_NLLS[1]-BETA_NLLS[2]*X[I]))**2
        *exp(-2*BETA_NLLS[1]-2*BETA_NLLS[2]*X[I])*(1/X[I])*(1|X[I])/nrow(Y);
end;
VAR_NLLS=inv(A_NLLS)*B_NLLS*inv(A_NLLS)/nrow(Y);
STE_NLLS=sqrt(vecdiag(VAR_NLLS));
/*****
3.4. Compute Weighted NLS estimates (by maximizing Q with QN method)
*****/
start QWNLLS(BETA) global(X,Y,BETA_NLLS);
    QWNLLS=0;
    do I=1 to nrow(Y);
        QWNLLS=QWNLLS-(Y[I]-exp(-BETA[1]-BETA[2]*X[I]))**2
            / (2*nrow(Y)*exp(-BETA_NLLS[1]-BETA_NLLS[2]*X[I]));
    end;
    return(QWNLLS);
finish QWNLLS;
call nlpqn(RC3,BETA_WNLLS,"QWNLLS",BETA_NLLS,{1,0});
BETA_WNLLS=BETA_WNLLS`;
U_WNLLS=j(nrow(Y),1.);
do I=1 to nrow(Y);
    U_WNLLS[I]=Y[I]-exp(-BETA_WNLLS[1]-BETA_WNLLS[2]*X[I]);
end;
U2_WNLLS=U_WNLLS##2;
A_WNLLS=0;
B_WNLLS=0;
do I=1 to nrow(Y);
    A_WNLLS=A_WNLLS+exp(-2*BETA_WNLLS[1]-2*BETA_WNLLS[2]*X[I])*(1/X[I])*(1|X[I])
        /exp(-BETA_WNLLS[1]-BETA_WNLLS[2]*X[I]);
    B_WNLLS=B_WNLLS+exp(-2*BETA_WNLLS[1]-2*BETA_WNLLS[2]*X[I])*(1/X[I])*(1|X[I])
        *U2_WNLLS[I]/exp(2*(-BETA_WNLLS[1]-BETA_WNLLS[2]*X[I]));
end;
VAR_WNLLS=inv(A_WNLLS)*B_WNLLS*inv(A_WNLLS);
STE_WNLLS=sqrt(vecdiag(VAR_WNLLS));
/*****
3.5. Compute Feasible Generalized NLS estimates (by maximizing Q with QN method)
*****/
start QFGNLLS(BETA) global(X,Y,BETA_NLLS);
    QFGNLLS=0;
    do I=1 to nrow(Y);
        QFGNLLS=QFGNLLS-(Y[I]-exp(-BETA[1]-BETA[2]*X[I]))**2
            / (2*nrow(Y)*exp(-2*BETA_NLLS[1]-2*BETA_NLLS[2]*X[I]));
    end;
    return(QFGNLLS);
finish QFGNLLS;
call nlpqn(RC3,BETA_FGNLLS,"QFGNLLS",BETA_NLLS,{1,0});
BETA_FGNLLS=BETA_FGNLLS`;
A_FGNLLS=0;
do I=1 to nrow(Y);

```



```

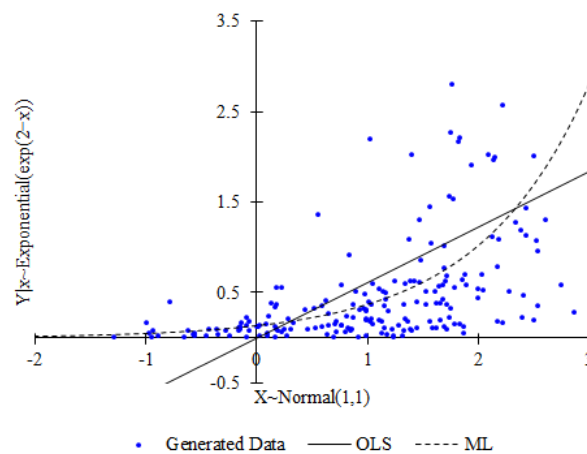
A_FGNLLS=A_FGNLLS+exp(-2*BETA_FGNLLS[1]-2*BETA_FGNLLS[2]*X[I])*(1/X[I])*(1|X[I])
/exp(-2*BETA_FGNLLS[1]-2*BETA_FGNLLS[2]*X[I]);
end;
VAR_FGNLLS=inv(A_FGNLLS);
STE_FGNLLS=sqrt(vecdiag(VAR_FGNLLS));
start RSQ_NL(BETA) global(X,Y);
U=j(nrow(Y),1,.);
do I=1 to nrow(Y);
U[I]=Y[I]-exp(-BETA[1]-BETA[2]*X[I]);
end;
RSQ_NL=1-U`*U/((Y-mean(Y))`*(Y-mean(Y)));
return(RSQ_NL);
finish RSQ_NL;
*print BETA_OLS[format=12.4] STE_OLS[format=12.4],
BETA_MLE[format=12.4] STE_MLE[format=12.4],
BETA_NLLS[format=12.4] STE_NLLS[format=12.4],
BETA_WNLLS[format=12.4] STE_WNLLS[format=12.4],
BETA_FGNLLS[format=12.4] STE_FGNLLS[format=12.4];
ESTIMATOR=(BETA_OLS[1]|BETA_MLE[1]|BETA_NLLS[1]|BETA_WNLLS[1]|BETA_FGNLLS[1])//
(STE_OLS[1]|STE_MLE[1]|STE_NLLS[1]|STE_WNLLS[1]|STE_FGNLLS[1])//
(BETA_OLS[2]|BETA_MLE[2]|BETA_NLLS[2]|BETA_WNLLS[2]|BETA_FGNLLS[2])//
(STE_OLS[2]|STE_MLE[2]|STE_NLLS[2]|STE_WNLLS[2]|STE_FGNLLS[2])//
(.|LNHAT(BETA_MLE)|LNHAT(BETA_NLLS)|LNHAT(BETA_WNLLS)|LNHAT(BETA_FGNLLS))//
(RSQ_OLS|RSQ_NL(BETA_MLE)|RSQ_NL(BETA_NLLS)|RSQ_NL(BETA_WNLLS)|RSQ_NL(BETA_FGNLLS));
print ESTIMATOR[colname={"OLS" "ML" "NLS" "WNLS" "FGNLS"}
rowname={"Constant" "x" "lnL" "R^2"} format=12.4];
quit;

```

SAS Result for the Simulation

The SAS System 17:04 Tuesday, October 25, 2016 614					
	OLS	ESTIMATOR ML	NLS	WNLS	FGNLS
Constant	-0.0025	1.9997	1.7780	1.9493	1.9968
x	0.0149	0.0143	0.1125	0.0289	0.0143
	0.6139	-1.0108	-0.9115	-0.9806	-1.0082
	0.0106	0.0101	0.0511	0.0187	0.0101
lnL	.	-190.9147	-304.6058	-197.2558	-190.9480
R^2	0.2526	0.3704	0.3802	0.3765	0.3712

Scatter Plot



Alarm: Midterm November 3, 2016

- Cover many things from nonlinear estimation
- Also cover sketches of proofs of consistency, Normality...

Numerical Methods of Optimization

Nonlinear estimators are usually only implicitly defined.

e.g.

$$\hat{\theta}_n \equiv \operatorname{argmax}_{\theta \in \Theta} \hat{Q}_n(\theta), \quad (\text{M Estimation})$$

$$\left. \frac{\partial \hat{Q}_n(\theta)}{\partial \theta^\top} \right|_{\theta = \hat{\theta}_n} = \mathbf{0}, \quad (\text{Z Estimation})$$

Important to know how computers compute $\hat{\theta}_n$: essentially two possibilities

1. Grid Search

- Pick many possible values of θ in some set
- Compute $\hat{Q}_n(\theta)$ for each values
- Pick $\hat{\theta}_n$ to be the value of θ that maximizes $\hat{Q}_n(\theta)$ /makes $\frac{\partial \hat{Q}_n(\theta)}{\partial \theta^\top}$ close to $\mathbf{0}$ on the grid
- This method works if the grid is adequate
- This method is very inefficient if θ is event of moderately high dimensional

e.g. Suppose $\theta \in \mathbb{R}^{10}$. Pick a grid of just 10 points (this grid is likely to be too sparse). Then $\hat{Q}_n(\theta)$ or $\frac{\partial \hat{Q}_n(\theta)}{\partial \theta^\top}$ will need to be evaluated 10^{10} =10 billion times

2. Iterative Methods

- Much more popular than grid search
- Suppose we want to estimate

$$\theta_0 \equiv \operatorname{argmax}_{\theta \in \Theta} E[\hat{Q}_n(\theta)], \quad \text{via } \hat{\theta}_n \equiv \operatorname{argmax}_{\theta \in \Theta} \hat{Q}_n(\theta)$$

- The idea is to update estimates of θ_0 via an algorithm
- Suppose $\hat{\theta}_s$ estimates θ_0 at iteration s
- The idea is to use $\hat{\theta}_s$ to generate another estimate $\hat{\theta}_{s+1}$ such that $\hat{Q}_n(\hat{\theta}_{s+1}) > \hat{Q}_n(\hat{\theta}_s)$
- Most iterative methods are gradient methods, i.e. we have for each iteration $s \in \{1, \dots, S\}$

$$\hat{\theta}_{s+1} = \hat{\theta}_s + \mathbf{A}_s \mathbf{g}_s$$

where \mathbf{A}_s is a matrix that depends on $\hat{\theta}_s$ and

$$\mathbf{g}_s = \left. \frac{\partial \hat{Q}_n(\theta)}{\partial \theta^\top} \right|_{\theta = \hat{\theta}_s}, \quad \text{i. e. the } \underline{\text{gradient}} \text{ vector at } \hat{\theta}_s$$

- Different gradient methods are distinguished by different choice of weighting matrix \mathbf{A}_s
e.g. the Newton–Raphson method uses

$$\mathbf{A}_s = -\mathbf{H}_s^{-1}, \quad \text{where } \mathbf{H}_s = \left. \frac{\partial^2 \hat{Q}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_s}$$

- In general, want \mathbf{A}_s to be positive definite with high probability. This makes $\hat{Q}_n(\hat{\boldsymbol{\theta}}_{s+1}) > \hat{Q}_n(\hat{\boldsymbol{\theta}}_s)$ with high probability
- To see this, $\hat{Q}_n(\hat{\boldsymbol{\theta}}_{s+1}) = \hat{Q}_n(\hat{\boldsymbol{\theta}}_s) + \mathbf{g}_s^\top \mathbf{A}_s \mathbf{g}_s + R_n$, where $R_n = o_p(\mathbf{g}_s^\top \mathbf{A}_s \mathbf{g}_s)$

Note: If $\|\mathbf{A}_s\|$ is too small, then the algorithm might a slow rate of convergence. If $\|\mathbf{A}_s\|$ is too large, then the algorithm might cause $\{\hat{\boldsymbol{\theta}}_s\}$ to overshoot the actual maximizer of $\hat{Q}_n(\boldsymbol{\theta})$

- Gradient methods often use step-size adjustments to reduce the risk of overshooting/undershooting

Take

$$\hat{\boldsymbol{\theta}}_{s+1} = \hat{\boldsymbol{\theta}}_s + \hat{\lambda}_s \mathbf{A}_s \mathbf{g}_s$$

where $\hat{\lambda}_s$ is a scalar step size chosen to maximize $\hat{Q}_n(\hat{\boldsymbol{\theta}}_{s+1})$

- This works by calculating $\mathbf{A}_s \mathbf{g}_s$, followed by calculating $\hat{Q}_n(\hat{\boldsymbol{\theta}}_s + \lambda \mathbf{A}_s \mathbf{g}_s) \forall$ values of λ in some chosen grid Λ . The final choice of $\hat{\lambda}_s$ is the true value $\lambda \in \Lambda$ that maximizes $\hat{Q}_n(\hat{\boldsymbol{\theta}}_s + \lambda \mathbf{A}_s \mathbf{g}_s)$
- This procedure minimizes the total number of iterations at the cost of doing grid searches between iterations
- Gradient methods are tend to converge to local maxima close to the chosen starting value $\hat{\boldsymbol{\theta}}_1$. If $\hat{Q}_n(\boldsymbol{\theta})$ is suspected to have local maxima, then one should try different starting values spaced far apart
- Iterations continued until (in the ideal case) converge, which is defined to occur when

(a) $|\hat{Q}_n(\hat{\boldsymbol{\theta}}_{s+1}) - \hat{Q}_n(\hat{\boldsymbol{\theta}}_s)|$ is small

(b) $\|\mathbf{A}_{s+1} \mathbf{g}_{t+1} - \mathbf{A}_s \mathbf{g}_s\|$ is small

(c) $\|\hat{\boldsymbol{\theta}}_{s+1} - \hat{\boldsymbol{\theta}}_s\|$ is small

- Ideally want all three of these convergence criteria to be very small, say $< 10^{-6}$
- Gradient methods also typically involve a preset maximum number of iterations: if the number of iterations exceeds this maximum without converging, then the algorithm stops and is said not to have converged
- If the algorithm takes S iterations to converge, then $\hat{\boldsymbol{\theta}}_s$ is a local maximum of $\hat{Q}_n(\boldsymbol{\theta})$: no guarantee that $\hat{\boldsymbol{\theta}}_s$ is a global maximum unless $\hat{Q}_n(\boldsymbol{\theta})$ is globally concave
- Can reduce the total number of iterations by picking a starting value $\hat{\boldsymbol{\theta}}_1$, that is a consistent estimate of $\boldsymbol{\theta}_0$

e.g. Newton–Raphson (NR)

This is a gradient method appropriate for cases where $\hat{Q}_n(\boldsymbol{\theta})$ is globally concave

$$\mathbf{A}_s = -\mathbf{H}_s^{-1}, \quad \text{where } \mathbf{H}_s = \left. \frac{\partial^2 \hat{Q}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_s}$$

(Note: \mathbf{H}_s^{-1} is guaranteed to exist for all possible values of $\hat{\boldsymbol{\theta}}_s$ if $\hat{Q}_n(\boldsymbol{\theta})$ is globally concave)

The NR method is motivated via a Taylor expansion

$$\hat{Q}_n(\boldsymbol{\theta}) = \hat{Q}_n(\hat{\boldsymbol{\theta}}_s) + \mathbf{g}_s^\top (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s) + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s)^\top \mathbf{H}_s (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s) + \hat{R}_n$$

which if \hat{R}_n is small, yields the approximation

$$\hat{Q}_n(\boldsymbol{\theta}) \approx \hat{Q}_n^*(\boldsymbol{\theta}) \equiv \hat{Q}_n(\hat{\boldsymbol{\theta}}_s) + \mathbf{g}_s^\top (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s) + \frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s)^\top \mathbf{A}_s (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s)$$

We see that $\hat{Q}_n^*(\boldsymbol{\theta})$ is globally maximized by solving

$$\begin{aligned} \frac{\partial \hat{Q}_n^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} = \mathbf{0} &\Leftrightarrow \mathbf{g}_s + \mathbf{H}_s (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_s) = \mathbf{0} \\ &\Leftrightarrow \boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_s - \mathbf{H}_s^{-1} \mathbf{g}_s \end{aligned}$$

which implies that $\hat{\boldsymbol{\theta}}_{s+1} = \hat{\boldsymbol{\theta}}_s - \mathbf{H}_s^{-1} \mathbf{g}_s$ is a global maximum (since \mathbf{H}_s is negative definite) of $\hat{Q}_n^*(\boldsymbol{\theta})$.

Note:

$$\begin{aligned} \hat{Q}_n(\hat{\boldsymbol{\theta}}_{s+1}) - \hat{Q}_n(\hat{\boldsymbol{\theta}}_s) &= -\mathbf{g}_s^\top \mathbf{H}_s^{-1} \mathbf{g}_s + \frac{1}{2} (-\mathbf{H}_s^\top \mathbf{g}_s)^\top \mathbf{H}_s (-\mathbf{H}_s^{-1} \mathbf{g}_s) + \hat{R}_n \\ &= -\frac{1}{2} \mathbf{g}_s^\top \mathbf{H}_s^{-1} \mathbf{g}_s + \hat{R}_n \end{aligned}$$

So if \hat{R}_n is small, then $\hat{Q}_n(\hat{\boldsymbol{\theta}}_{s+1}) - \hat{Q}_n(\hat{\boldsymbol{\theta}}_s) \approx -\frac{1}{2} \mathbf{g}_s^\top \mathbf{H}_s^{-1} \mathbf{g}_s$, which is positive if \mathbf{H}_s is negative definite (exercise: prove this).

Note: If $\hat{\boldsymbol{\theta}}_s$ is a local maximum, then \mathbf{H}_s is negative semidefinite by definition.

- If $\hat{\boldsymbol{\theta}}_s$ is far from a local maximum, then it is unclear if \mathbf{H}_s is even a definite matrix without assumptions on the shape of $\hat{Q}_n(\boldsymbol{\theta})$: in such cases, NR may fail
- If \mathbf{H}_s is singular, then $\hat{\boldsymbol{\theta}}_{s+1}$ is not computable and algorithm fails
- For these reasons, NR works best when $\hat{Q}_n(\boldsymbol{\theta})$ is globally concave

e.g. Method of Scoring (MS)

This is a variant of NR, where instead of $\mathbf{A}_s = -\mathbf{H}_s^{-1}$, we use $\mathbf{A}_s = -\bar{\mathbf{H}}_s^{-1}$, where

$$\bar{\mathbf{H}}_s = E \left[\left. \frac{\partial^2 \hat{Q}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_s} \right], \quad \text{or a consistent estimate of } E \left[\left. \frac{\partial^2 \hat{Q}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right]$$

MS is often applied to compute MLEs, i.e. cases where

$$\hat{Q}_n(\boldsymbol{\theta}) = \frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta})$$

Recall by the information identity,

$$-E \left[\frac{\partial^2 \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] = E \left[\frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right]$$

which is positive definite in view of the score identity.

Therefore, $\hat{\mathbf{H}}_s$ is likely to be negative definite (for MLE computation) when $\hat{\boldsymbol{\theta}}_s$ is close to $\boldsymbol{\theta}_0$ (e.g. $\hat{\boldsymbol{\theta}}_s$ is consistent).

e.g. BHHH

Another variant of NR, due to Berndt, Hall, Hall and Hausman (1974): Suppose

$$\hat{Q}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n q(\mathbf{Z}_i, \boldsymbol{\theta})$$

where $q(\mathbf{z}, \boldsymbol{\theta})$ to be continuously differentiable in $\boldsymbol{\theta}$ for all \mathbf{z} . The BHHH algorithm uses $\mathbf{A}_s = -\hat{\mathbf{H}}_s^{-1}$, where

$$\hat{\mathbf{H}}_s = -\frac{1}{n} \sum_{i=1}^n \frac{\partial q(\mathbf{z}_i, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial q(\mathbf{z}_i, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_s}$$

(i.e. no second-order derivative)

BHHH, like MS, is usually applied to compute MLEs: Suppose

$$\hat{Q}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log f(\mathbf{Z}_i, \boldsymbol{\theta}), \quad \text{where } \mathbf{Z}_1, \dots, \mathbf{Z}_n \text{ are iid with density } f(\cdot; \boldsymbol{\theta})$$

By the information identity, we have

$$\begin{aligned} \frac{1}{n} E \left[\frac{\partial^2 \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] &= -\frac{1}{n} E \left[\frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial \log \hat{L}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] \\ &= -E \left[\frac{\partial \log f(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \frac{\partial \log f(\mathbf{Z}_1, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right] \end{aligned}$$

which is negative definite by score identity.

Alarm: Midterm November 3, 2016

- Cover many things from nonlinear estimation
 - Especially, one should know how to derive the asymptotic distribution of $\sqrt{n}(\hat{\theta}_n - \theta_0)$ or the sandwich-style variance $A_0^{-1}B_0A_0^{-1}$, etc.
- Also cover sketches of proofs of consistency, Normality...

Gradient methods: Goal is to compute

$$\hat{\theta}_n \equiv \underset{\theta}{\operatorname{argmax}} \hat{Q}_n(\theta)$$

by means of an algorithm that generates a sequence

$$\hat{\theta}_s \rightarrow \hat{\theta}_n$$

where for each $s \in \{1, \dots, S\}$

$$\hat{\theta}_{s+1} = \hat{\theta}_s + \mathbf{A}_s \mathbf{g}_s$$

e.g. c) Steepest ascent: The simplest of all gradient methods

Set $\mathbf{A}_s = \mathbf{I}_d$ and $\hat{\theta}_{s+1} = \hat{\theta}_s + \hat{\lambda}_s \mathbf{g}_s$ where the step size $\hat{\lambda}_s$ maximizes $\hat{Q}_n(\hat{\theta}_s + \lambda \mathbf{g}_s)$ over λ in some grid Λ . Usually,

$$\hat{\lambda}_s = \frac{\mathbf{g}_s^\top \mathbf{g}_s}{\mathbf{g}_s^\top \mathbf{H}_s \mathbf{g}_s}, \quad \text{where } \mathbf{H}_s = \left. \frac{\partial^2 \hat{Q}_n(\theta)}{\partial \theta \partial \theta^\top} \right|_{\theta = \hat{\theta}_s}$$

(for this setting of $\hat{\lambda}_s$ to work, need $\hat{\lambda}_s \mathbf{I}_d$ to be negative definite, i.e. $\hat{\lambda}_s < 0$, i.e. \mathbf{H}_s to be negative definite)

Review: Why does \mathbf{A}_s need to be positive definite?

$$\begin{aligned} \hat{Q}_n(\hat{\theta}_{s+1}) - \hat{Q}_n(\hat{\theta}_s) &= \mathbf{g}_s^\top \mathbf{A}_s \mathbf{g}_s + \hat{R}_n \\ &\approx \mathbf{g}_s^\top \mathbf{A}_s \mathbf{g}_s, \quad \text{if } \hat{R}_n \text{ is small} \end{aligned}$$

Need $\mathbf{g}_s^\top \mathbf{A}_s \mathbf{g}_s > 0$, i.e. \mathbf{A}_s should be positive definite.

e.g. d) Gauss–Newton (GN): Gradient method applicable to NLLS

Suppose $E[Y|X]=g(\mathbf{X},\boldsymbol{\beta})$. Observe $[\mathbf{X}_i^\top Y_i]^\top$ ($i=1,\dots,n$) iid. The GN method sets $\widehat{\boldsymbol{\beta}}_{s+1} - \widehat{\boldsymbol{\beta}}_s$ equal to the OLS estimates in a regression given by

$$Y_i - g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s) = \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} \boldsymbol{\beta} + v_i$$

i.e. $\widehat{\boldsymbol{\beta}}_{s+1}$ is the OLS estimate for a regression given by

$$Y_i - g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s) - \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} \widehat{\boldsymbol{\beta}}_s = \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} \boldsymbol{\beta} + v_i$$

GN is motivated by the following Taylor expansion

$$g(\mathbf{X}_i, \boldsymbol{\beta}) = g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s) + \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_s) + R_n$$

Assume that R_n is small, substitute the approximation

$$g^*(\mathbf{X}_i, \boldsymbol{\beta}) \equiv g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s) + \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_s)$$

into the NLLS criterion function to get

$$\widehat{Q}_n(\boldsymbol{\beta}) \approx -\frac{1}{2n} \sum_{i=1}^n \left[Y_i - g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s) - \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_s) \right]^2$$

which is equivalent to the OLS criterion function for a regression of $Y_i - g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s)$ on $\left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s}$ with parameter vector $\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_s$. This implies that

$$\widehat{\boldsymbol{\beta}}_{s+1} = \widehat{\boldsymbol{\beta}}_s + \left(\sum_{i=1}^n \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} \right)^{-1} \sum_{i=1}^n \left. \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \right|_{\boldsymbol{\beta}=\widehat{\boldsymbol{\beta}}_s} (Y_i - g(\mathbf{X}_i, \widehat{\boldsymbol{\beta}}_s))$$

This shows that GN is a special case of MS applied to NLLS.

$$\begin{aligned} \text{For } \widehat{Q}_n(\boldsymbol{\beta}) &= -\frac{1}{2n} \sum_{i=1}^n (Y_i - g(\mathbf{X}_i, \boldsymbol{\beta}))^2, \quad \text{we have} \\ \frac{\partial \widehat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} &= \frac{1}{n} \sum_{i=1}^n (Y_i - g(\mathbf{X}_i, \boldsymbol{\beta})) \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \\ \frac{\partial^2 \widehat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} &= -\frac{1}{n} \sum_{i=1}^n \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} + \frac{1}{n} \sum_{i=1}^n (Y_i - g(\mathbf{X}_i, \boldsymbol{\beta})) \frac{\partial^2 g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \\ \Rightarrow E \left[\frac{\partial^2 \widehat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \right] &= -E \left[\frac{\partial g(\mathbf{X}_1, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \frac{\partial g(\mathbf{X}_1, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right] \end{aligned}$$

which is approximable as

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} \bigg|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_s} \approx \frac{\partial g(\mathbf{X}_i, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \bigg|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}_s}$$

(compare with the expression for $\hat{\boldsymbol{\beta}}_{s+1}$)

Example: Suppose we observe $[\mathbf{X}_i^\top \ Y_i]^\top$ ($i=1, \dots, n$) iid from a model where

$$E[Y_i | \mathbf{X}_i] = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0}, \quad \text{where } \boldsymbol{\beta}_0 \in \mathbb{R}^d$$

The NLLS estimator of $\boldsymbol{\beta}_0$ is

$$\hat{\boldsymbol{\beta}}_{\text{NLLS}} \equiv \underset{\boldsymbol{\beta} \in \mathbb{R}^d}{\operatorname{argmax}} \hat{Q}_n(\boldsymbol{\beta}), \quad \text{where } \hat{Q}_n(\boldsymbol{\beta}) = -\frac{1}{2n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}})^2$$

We have

$$\begin{aligned} \mathbf{g}(\boldsymbol{\beta}) &\equiv \frac{\partial \hat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^\top} = \frac{1}{n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}}) e^{\mathbf{X}_i^\top \boldsymbol{\beta}} \mathbf{X}_i \\ \mathbf{H}(\boldsymbol{\beta}) &\equiv \frac{\partial^2 \hat{Q}_n(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} = -\frac{1}{n} \sum_{i=1}^n e^{2\mathbf{X}_i^\top \boldsymbol{\beta}} \mathbf{X}_i \mathbf{X}_i^\top + \frac{1}{n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \boldsymbol{\beta}}) e^{\mathbf{X}_i^\top \boldsymbol{\beta}} \mathbf{X}_i \mathbf{X}_i^\top \end{aligned}$$

NR sets

$$\hat{\boldsymbol{\beta}}_{s+1} = \hat{\boldsymbol{\beta}}_s - \mathbf{H}^{-1}(\hat{\boldsymbol{\beta}}_s) \mathbf{g}(\hat{\boldsymbol{\beta}}_s)$$

MS exploits the fact that

$$E_{\boldsymbol{\beta}_0}[\mathbf{H}(\boldsymbol{\beta}_0)] = -E_{\boldsymbol{\beta}_0} \left[e^{2\mathbf{X}_i^\top \boldsymbol{\beta}_0} \mathbf{X}_i \mathbf{X}_i^\top \right]$$

which can be approximated by

$$\bar{\mathbf{H}}(\hat{\boldsymbol{\beta}}_s) \equiv -\frac{1}{n} \sum_{i=1}^n e^{2\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_s} \mathbf{X}_i \mathbf{X}_i^\top$$

Therefore MS sets

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{s+1} &= \hat{\boldsymbol{\beta}}_s - \bar{\mathbf{H}}^{-1}(\hat{\boldsymbol{\beta}}_s) \mathbf{g}(\hat{\boldsymbol{\beta}}_s) \\ &= \hat{\boldsymbol{\beta}}_s + \left(\frac{1}{n} \sum_{i=1}^n e^{2\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_s} \mathbf{X}_i \mathbf{X}_i^\top \right)^{-1} \frac{1}{n} \sum_{i=1}^n (Y_i - e^{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_s}) e^{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_s} \mathbf{X}_i \end{aligned}$$

i.e. $\hat{\boldsymbol{\beta}}_{s+1} - \hat{\boldsymbol{\beta}}_s$ is the OLS estimate for a regression of $Y_i - e^{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_s}$ on $e^{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_s} \mathbf{X}_i$. This shows that in this case, MS is GN.

3. Non-Gradient Methods

a) EM algorithm: The expectation maximization (EM) algorithm is an iterative method for computing MLEs (Hartley, 1958, Biometrics).

In Econometrics, the EM algorithm is usually associated with MLEs of latent-variable models, i.e. a model where a matrix \mathbf{Z} of observable variables is related to a vector \mathbf{y}^* of unobservable variables in such a way that \mathbf{y}^* uniquely determines \mathbf{Z} , but \mathbf{Z} does not uniquely define \mathbf{y}^* .

e.g. Suppose we observe $[\mathbf{D}_i \mathbf{X}_i^\top \mathbf{Y}_i]^\top$ ($i=1, \dots, n$) where for each i

$$\mathbf{Y}_i^* = \mathbf{X}_i^\top \boldsymbol{\beta}_0 + u_i, \quad \mathbf{D}_i = 1\{\mathbf{Y}_i^* > 0\}, \quad \mathbf{Y}_i = \mathbf{D}_i \mathbf{Y}_i^*$$

(and where u_1, \dots, u_n are iid $N(0, \sigma_0^2)$ and unobservable)

This is the standard Tobit model. In this case,

$$\mathbf{Y}^* = \begin{bmatrix} \mathbf{Y}_1^* \\ \vdots \\ \mathbf{Y}_n^* \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} \mathbf{D}_1 & \mathbf{X}_1^\top & \mathbf{Y}_1 \\ \vdots & \vdots & \vdots \\ \mathbf{D}_n & \mathbf{X}_n^\top & \mathbf{Y}_n \end{bmatrix}$$

Note that observations for which $\mathbf{D}_i=0$ do not allow one to infer anything about \mathbf{Y}_i^* apart from the fact that $\mathbf{Y}_i^* \leq 0$. In general, suppose \mathbf{Y}^* (i.e. the latent variable) has joint density $f(\mathbf{y}^*)$. Let \mathbf{Z} have density $g(\mathbf{z})$.

$$\text{Let } k(\mathbf{y}^*|\mathbf{z}) \equiv \frac{f(\mathbf{y}^*)}{g(\mathbf{z})}$$

The joint density of \mathbf{Y}^* and \mathbf{Z} is

$$f(\mathbf{y}^*, \mathbf{z}) = f(\mathbf{z}|\mathbf{y}^*)f(\mathbf{y}^*) = f(\mathbf{y}^*)$$

because $f(\mathbf{z}|\mathbf{y}^*) \equiv 1$ in view of the fact that \mathbf{y}^* uniquely determines \mathbf{Z} .

Assume that $f(\cdot)$, $g(\cdot)$ and $k(\cdot)$ all depend on a parameter $\boldsymbol{\theta} \in \mathbb{R}^d$, i.e.

$$f(\mathbf{y}^*) = f(\mathbf{y}^*; \boldsymbol{\theta}), \quad g(\mathbf{z}) = g(\mathbf{z}; \boldsymbol{\theta}), \quad k(\mathbf{y}^*|\mathbf{z}) = k(\mathbf{y}^*|\mathbf{z}; \boldsymbol{\theta})$$

Want to compute the MLE, i.e.

$$\hat{\boldsymbol{\theta}}_{\text{MLE}} \equiv \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta})$$

where $\hat{L}_n(\boldsymbol{\theta}) = g(\mathbf{z}; \boldsymbol{\theta}) = \frac{f(\mathbf{y}^*; \boldsymbol{\theta})}{k(\mathbf{y}^*|\mathbf{z}; \boldsymbol{\theta})}$. So

$$\frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta}) = \frac{1}{n} \log f(\mathbf{y}^*; \boldsymbol{\theta}) - \frac{1}{n} \log k(\mathbf{y}^*|\mathbf{z}; \boldsymbol{\theta})$$

$$\text{Let } \bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s) \equiv E_{\boldsymbol{\theta}_s} \left[\frac{1}{n} \log f(\mathbf{y}^*; \boldsymbol{\theta}) \mid \mathbf{Z} \right]$$

The EM algorithm maximizes $\frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta})$ by maximizing $\bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)$ with respect to $\boldsymbol{\theta}$.

Let

$$\bar{S}(\boldsymbol{\theta}; \boldsymbol{\theta}_s) \equiv E_{\boldsymbol{\theta}_s} \left[\frac{1}{n} \log k(\mathbf{y}^* | \mathbf{z}; \boldsymbol{\theta}) \mid \mathbf{Z} \right]$$

Note:

$$\bar{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}_s) \equiv E_{\boldsymbol{\theta}_s} \left[\frac{1}{n} \log g(\mathbf{z}; \boldsymbol{\theta}) \mid \mathbf{Z} \right] = \frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta}) = \hat{Q}_n(\boldsymbol{\theta})$$

So

$$\hat{Q}_n(\boldsymbol{\theta}) = \bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s) - \bar{S}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)$$

By Jensen's inequality

$$\begin{aligned} & \bar{S}(\boldsymbol{\theta}; \boldsymbol{\theta}_s) < \bar{S}(\boldsymbol{\theta}_s; \boldsymbol{\theta}_s), \quad \text{if } \boldsymbol{\theta} \neq \boldsymbol{\theta}_s, \quad \dots (*) \\ & \left(\text{notice: } E_{\boldsymbol{\theta}_0} \left[\log \frac{f(\mathbf{z}; \boldsymbol{\theta})}{f(\mathbf{z}; \boldsymbol{\theta}_0)} \right] < \log E_{\boldsymbol{\theta}_0} \left[\frac{f(\mathbf{z}; \boldsymbol{\theta})}{f(\mathbf{z}; \boldsymbol{\theta}_0)} \right] = 0, \quad \forall \boldsymbol{\theta} \neq \boldsymbol{\theta}_0 \right) \end{aligned}$$

For a given $\boldsymbol{\theta}_s$, let $\boldsymbol{\theta}_{s+1}^*$ maximizes $\bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)$

Then

$$\begin{aligned} & \hat{Q}_n(\boldsymbol{\theta}_{s+1}^*) = \bar{R}(\boldsymbol{\theta}_{s+1}^*; \boldsymbol{\theta}_s) - \bar{S}(\boldsymbol{\theta}_{s+1}^*; \boldsymbol{\theta}_s) \\ & \text{but } \bar{R}(\boldsymbol{\theta}_{s+1}^*; \boldsymbol{\theta}_s) \geq \bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s), \quad \forall \boldsymbol{\theta} \\ & \text{while } \bar{S}(\boldsymbol{\theta}_{s+1}^*; \boldsymbol{\theta}_s) \leq \bar{S}(\boldsymbol{\theta}_s; \boldsymbol{\theta}_s), \quad \text{by } (*) \\ & \text{hence } \hat{Q}_n(\boldsymbol{\theta}_{s+1}^*) \geq \hat{Q}_n(\boldsymbol{\theta}_s) \end{aligned}$$

This shows that the log-likelihood never decreases under the EM algorithm. In addition, if $\hat{L}_n(\boldsymbol{\theta})$ is bounded, then $\lim_{s \rightarrow \infty} \frac{1}{n} \log \hat{L}_n(\hat{\boldsymbol{\theta}}_s)$ exists.

To see this, suppose $\boldsymbol{\theta}^*$ is a stationary point of the EM algorithm. Then $\boldsymbol{\theta}^*$ is a stationary point of $\hat{Q}_n(\boldsymbol{\theta})$. Suppose that $\hat{Q}_n(\boldsymbol{\theta})$ is differentiable. Then

$$\left. \frac{\partial \hat{Q}_n(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_s} = \left. \frac{\partial \bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_s} - \left. \frac{\partial \bar{S}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_s}$$

But

$$\left. \frac{\partial \bar{S}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)}{\partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_s} = \mathbf{0}, \quad (\text{because } \bar{S}(\boldsymbol{\theta}; \boldsymbol{\theta}_s) < \bar{S}(\boldsymbol{\theta}_s; \boldsymbol{\theta}_s) \forall \boldsymbol{\theta} \in \boldsymbol{\theta}_s)$$

So if $\boldsymbol{\theta}_s$ is a stationary point of $\bar{R}(\boldsymbol{\theta}; \boldsymbol{\theta}_s)$ (i.e. a stationary point of the EM algorithm, it is also a stationary point of $\hat{Q}_n(\boldsymbol{\theta}) = \frac{1}{n} \log \hat{L}_n(\boldsymbol{\theta})$)

Latent Variable Models: Are characterized by a random matrix \underline{Z} of observables uniquely determined by another random matrix \underline{Z}^* of un-observables and such that \underline{Z} does not uniquely determine \underline{Z}^*

Notation: $c, \underline{c}, \underline{C}$ for scalar, column vector, matrix; $Y, \underline{Y}, \underline{Y}$ for random variable, column vector, matrix, respectively. Until now I have used expressions such as $c, \underline{c}, \underline{C}$ and $\tilde{y}, \tilde{\underline{y}}, \tilde{\underline{Y}}$, respectively.

e.g. Standard Tobit model

Observe $[D_i \underline{X}_i^\top Y_i]^\top$ ($i=1, \dots, n$) where for each i

$$Y_i^* = \underline{X}_i^\top \beta_0 + U_i, \quad Y_i = D_i Y_i^*, \quad D_i = 1\{Y_i^* > 0\}$$

and where U_1, \dots, U_n iid $N(0, \sigma_0^2)$. In this case

$$\underline{Z}^* = \begin{bmatrix} \underline{X}_1^\top & Y_1^* \\ \vdots & \vdots \\ \underline{X}_n^\top & Y_n^* \end{bmatrix}, \quad \underline{Z} = \begin{bmatrix} D_1 & \underline{X}_1^\top & Y_1 \\ \vdots & \vdots & \vdots \\ D_n & \underline{X}_n^\top & Y_n \end{bmatrix}$$

Note: Observations for which $D_i=0$ do not allow one to infer anything about the corresponding Y_i^* apart from the fact that $Y_i^* \leq 0$.

EM algorithm for computing the MLE of latent variable models

Let \underline{Y}^* (the latent variable) have joint density $f(\underline{y}^*)$. Let \underline{Z} (observed) have density $g(\underline{z})$.

Let $k(\underline{y}^* | \underline{z}) \equiv \frac{f(\underline{y}^*)}{g(\underline{z})}$. Note that the joint density of \underline{Y}^* and \underline{Z} is $f(\underline{y}^*, \underline{z}) = f(\underline{z} | \underline{y}^*) f(\underline{y}^*) = f(\underline{y}^*)$ since $f(\underline{z} | \underline{y}^*) \equiv 1$.

Assume that $f(\underline{y}^*)$, $g(\underline{z})$, $k(\underline{y}^* | \underline{z})$ all depend on a parameter $\underline{\theta}_0 \in \mathbb{R}^d$. Let $\underline{\theta} \in \mathbb{R}^d$ and write

$$f(\underline{y}^*) = f(\underline{y}^*; \underline{\theta}), \quad g(\underline{z}) = g(\underline{z}; \underline{\theta}), \quad k(\underline{y}^* | \underline{z}) = k(\underline{y}^* | \underline{z}; \underline{\theta})$$

Let the MLE be

$$\hat{\underline{\theta}}_{ML} \equiv \operatorname{argmax}_{\underline{\theta} \in \mathbb{R}^d} \frac{1}{n} \log \hat{L}_n(\underline{\theta})$$

where

$$\hat{L}_n(\underline{\theta}) = g(\underline{Z}; \underline{\theta}) = \frac{f(\underline{Y}^*; \underline{\theta})}{k(\underline{Y}^* | \underline{Z}; \underline{\theta})}, \quad \Rightarrow \frac{1}{n} \log \hat{L}_n(\underline{\theta}) = \frac{1}{n} \log f(\underline{Y}^*; \underline{\theta}) - \frac{1}{n} \log k(\underline{Y}^* | \underline{Z}; \underline{\theta})$$

Let

$$\bar{R}(\underline{\theta}; \underline{\theta}_s) \equiv E_{\underline{\theta}_s} \left[\frac{1}{n} \log f(\underline{Y}^*; \underline{\theta}) \mid \underline{Z} \right]$$

i.e. a conditional mean given \underline{z} where $\underline{\theta}_s$ is treated as the true value of $\underline{\theta}$.

The EM algorithm maximizes $\frac{1}{n} \log \hat{L}_n(\underline{\theta})$ by maximizing $\bar{R}(\underline{\theta}; \underline{\theta}_s)$ with respect to $\underline{\theta}$. To see why this works, let

$$\bar{S}(\underline{\theta}; \underline{\theta}_s) \equiv E_{\underline{\theta}_s} \left[\frac{1}{n} \log k(Y^* | \underline{Z}; \underline{\theta}) \mid \underline{Z} \right]$$

and observe that

$$\bar{Q}(\underline{\theta}; \underline{\theta}_s) \equiv E_{\underline{\theta}_s} \left[\frac{1}{n} \log g(\underline{Z}; \underline{\theta}) \mid \underline{Z} \right] = \frac{1}{n} \log \hat{L}_n(\underline{\theta})$$

So

$$\hat{Q}_n(\underline{\theta}) = \frac{1}{n} \log \hat{L}_n(\underline{\theta}) = \bar{R}(\underline{\theta}; \underline{\theta}_s) - \bar{S}(\underline{\theta}; \underline{\theta}_s)$$

By Jensen's inequality

$$\bar{S}(\underline{\theta}; \underline{\theta}_s) < \bar{S}(\underline{\theta}_s; \underline{\theta}_s), \quad \forall \underline{\theta} \neq \underline{\theta}_s, \quad \dots (*)$$

This is because

$$E_{\underline{\theta}_0} \left[\log \left(\frac{f(\underline{Z}; \underline{\theta})}{f(\underline{Z}; \underline{\theta}_0)} \right) \right] < \log E_{\underline{\theta}_0} \left[\frac{f(\underline{Z}; \underline{\theta})}{f(\underline{Z}; \underline{\theta}_0)} \right] = 0, \quad \text{when } \underline{\theta} \neq \underline{\theta}_0$$

since $\log E_{\underline{\theta}_0}[1] = \log 1 = 0$.

For a given $\underline{\theta}_s$, let $\underline{\theta}_{s+1}^*$ maximizes $\bar{R}(\underline{\theta}, \underline{\theta}_s)$. Then

$$\hat{Q}_n(\underline{\theta}_{s+1}^*) = \bar{R}(\underline{\theta}_{s+1}^*; \underline{\theta}_s) - \bar{S}(\underline{\theta}_{s+1}^*; \underline{\theta}_s)$$

But

$$\bar{R}(\underline{\theta}_{s+1}^*; \underline{\theta}_s) \geq \bar{R}(\underline{\theta}; \underline{\theta}_s), \quad \forall \underline{\theta} \in \mathbb{R}^d$$

while

$$\bar{S}(\underline{\theta}_{s+1}^*; \underline{\theta}_s) \leq \bar{S}(\underline{\theta}_s; \underline{\theta}_s), \quad \text{by } (*)$$

which implies that $\hat{Q}_n(\underline{\theta}_{s+1}^*) \geq \hat{Q}_n(\underline{\theta}_s)$

In addition, if $\hat{L}_n(\underline{\theta})$ is bounded, then $\lim_{s \rightarrow \infty} \frac{1}{n} \log \hat{L}_n(\underline{\theta}_s)$ exists. To see this, suppose $\underline{\theta}^*$ is a stationary point of the EM algorithm. Then $\underline{\theta}^*$ is a stationary point of $\hat{Q}_n(\underline{\theta})$.

If $\hat{Q}_n(\underline{\theta})$ is twice differentiable, then

$$\left. \frac{\partial \hat{Q}_n(\underline{\theta})}{\partial \underline{\theta}^T} \right|_{\underline{\theta}=\underline{\theta}_s} = \left. \frac{\partial \bar{R}(\underline{\theta}; \underline{\theta}_s)}{\partial \underline{\theta}^T} \right|_{\underline{\theta}=\underline{\theta}_s} - \left. \frac{\partial \bar{S}(\underline{\theta}; \underline{\theta}_s)}{\partial \underline{\theta}^T} \right|_{\underline{\theta}=\underline{\theta}_s}$$

But $\underline{\theta}_s$ is a global maximum of $\bar{S}(\underline{\theta}, \underline{\theta}_s)$, so

$$\left. \frac{\partial \bar{S}(\underline{\theta}; \underline{\theta}_s)}{\partial \underline{\theta}^T} \right|_{\underline{\theta}=\underline{\theta}_s} = \underline{0}$$

So if $\underline{\theta}_s$ is a stationary point of $\bar{R}(\underline{\theta}; \underline{\theta}_s)$, it is also a stationary point of $\hat{Q}_n(\underline{\theta})$.

e.g. Consider the standard Tobit model

Suppose we observe $[D_i \ X_i^\top \ Y_i]^\top$ ($i=1, \dots, n$) where for each i

$$\begin{aligned} Y_i^* &= X_i^\top \beta_0 + U_i \\ D_i &= 1\{Y_i^* > 0\} \\ Y_i &= Y_i^* D_i \end{aligned}$$

where U_1, \dots, U_n are iid $N(0, \sigma_0^2)$ and unobserved. Let

$$\begin{aligned} \underline{\theta} &\equiv [\underline{\beta}^\top \ \sigma^2]^\top \\ \underline{\theta}_0 &\equiv [\underline{\beta}_0^\top \ \sigma^2]^\top \\ \underline{Y}^* &\equiv [Y_1^* \ \dots \ Y_n^*]^\top \\ \underline{Z} &= \begin{bmatrix} D_1 & X_1^\top & Y_1 \\ \vdots & \vdots & \vdots \\ D_n & X_n^\top & Y_n \end{bmatrix} \end{aligned}$$

We have

$$\begin{aligned} \log f(\underline{Y}^*; \underline{\theta}) &= -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i^* - X_i^\top \underline{\beta})^2 \\ E_{\underline{\theta}_s} \left[\frac{1}{n} \log f(\underline{Y}^*; \underline{\theta}) \mid \underline{Z} \right] &= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2 n} \sum_{\{i: D_i=1\}} (Y_i - X_i^\top \underline{\beta})^2 \\ &\quad - \frac{1}{2\sigma^2 n} \sum_{\{i: D_i=0\}} E_{\underline{\theta}_s} [(Y_i^* - X_i^\top \underline{\beta})^2 \mid D_i = 0, X_i] \\ &= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 \\ &\quad - \frac{1}{2n\sigma^2} \sum_{\{i: D_i=1\}} (Y_i - X_i^\top \underline{\beta})^2 - \frac{1}{2\sigma^2 n} \sum_{\{i: D_i=0\}} \text{Var}_{\underline{\theta}_s} [Y_i^* \mid D_i = 0, X_i] \\ &\quad - \frac{1}{2\sigma^2 n} \sum_{\{i: D_i=0\}} (E_{\underline{\theta}_s} [Y_i^* \mid D_i = 0, X_i] - X_i^\top \underline{\beta})^2, \quad \dots (**) \end{aligned}$$

where

$$\begin{aligned} E_{\underline{\theta}_s} [Y_i^* \mid D_i = 0, X_i] &= X_i^\top \underline{\beta}_s - \frac{\sigma_s \phi_{is}}{1 - \Phi_{is}} \equiv \bar{Y}_{i0} \\ \text{Var}_{\underline{\theta}_s} [Y_i^* \mid D_i = 0, X_i] &= \sigma_s^2 + X_i^\top \underline{\beta}_s \frac{\sigma_s \phi_{is}}{1 - \Phi_{is}} - \left(\frac{\sigma_s \phi_{is}}{1 - \Phi_{is}} \right)^2 \end{aligned}$$

where

$$\phi_{is} = \phi \left(\frac{X_i^\top \underline{\beta}_s}{\sigma_s} \right), \quad \Phi_{is} = \Phi \left(\frac{X_i^\top \underline{\beta}_s}{\sigma_s} \right)$$

(e.g. Hayashi (2000), §8.2)

Since $Y_i^* | X_i \sim N(\underline{X}_i^\top \underline{\beta}_s, \sigma_s^2)$ when computing $E_{\underline{\theta}_s}[Y_i^* | D_i = 0, \underline{X}_i]$ and $\text{Var}_{\underline{\theta}_s}[Y_i^* | D_i = 0, \underline{X}_i]$, from (**), $\underline{\beta}_{s+1}$ can be computed as follows.

Let \underline{Y} be the vector consisting of Y_1, \dots, Y_n ; assumed all positive WOLOG.

Let $\underline{Y}_0 \equiv [\underline{Y}_{i0}]$ consist of $n - n_1$ elements.

Let $\underline{X} \equiv [\underline{X}_1 \dots \underline{X}_n]^\top$ as usual.

Then by (**),

$$E_{\underline{\theta}_s} \left[\frac{1}{n} \log f(\underline{Y}^*; \underline{\theta}) \mid \underline{Z} \right]$$

is maximized w.r.t. $\underline{\beta}$ by OLS, i.e.

$$\underline{\beta}_{s+1} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \left[\frac{\underline{Y}}{\underline{Y}_0} \right], \quad \dots (\dagger)$$

Similarly, by (**),

$$E_{\underline{\theta}_s} \left[\frac{1}{n} \log f(\underline{Y}^*; \underline{\theta}) \mid \underline{Z} \right]$$

is maximized w.r.t. σ^2 by OLS i.e.

$$\sigma_{s+1}^2 = \frac{1}{n} \left[\sum_{\{i: D_i=1\}} (Y_i - \underline{X}_i^\top \underline{\beta}_{s+1})^2 + \sum_{\{i: D_i=0\}} (\bar{Y}_{i0} - \underline{X}_i^\top \underline{\beta}_{s+1})^2 + \sum_{\{i: D_i=0\}} \text{Var}_{\underline{\theta}_s}[Y_i^* | D_i = 0, \underline{X}_i] \right], \quad \dots (\ddagger)$$

We can show that $\hat{\underline{\theta}}_{ML}$ is the equilibrium solution of the iterative progress given by (\dagger) , (\ddagger) . Partition

$$\underline{X} = \begin{bmatrix} \underline{X}_1 \\ \underline{X}_0 \end{bmatrix}$$

where \underline{X}_1 is $(n_1 \times d)$ and \underline{X}_0 is $[(n - n_1) \times d]$. Insert $\hat{\underline{\theta}}_{ML} \equiv [\hat{\underline{\beta}}_{ML}^\top \hat{\sigma}_{ML}^2]^\top$ into both sides of (\dagger) to get

$$\underline{X}_1^\top \underline{X}_1 \hat{\underline{\beta}}_{ML} \equiv \underline{X}_1^\top \underline{Y} - \underline{X}_0^\top \hat{\underline{Y}}_{ML} \Leftrightarrow -\underline{X}_0^\top \hat{\underline{Y}}_{ML} = -\underline{X}_1^\top (\underline{Y} - \underline{X}_1 \hat{\underline{\beta}}_{ML}), \quad \dots (***)$$

Where $\hat{\underline{Y}}_{ML}$ is an $[(n - n_1) \times 1]$ -vector whose j -th element is

$$\frac{\hat{\sigma}_{ML} \Phi(\underline{X}_j^\top \hat{\underline{\beta}}_{ML} / \hat{\sigma}_{ML})}{1 - \Phi(\underline{X}_j^\top \hat{\underline{\beta}}_{ML} / \hat{\sigma}_{ML})}$$

We know that if $\left. \frac{\partial}{\partial \underline{\beta}^\top} \frac{1}{n} \log \hat{L}_n(\underline{\theta}) \right|_{\underline{\theta} = \hat{\underline{\theta}}_{ML}} = \underline{0}$, so

$$-\sigma \sum_{\{i: D_i=0\}} \frac{\Phi_i}{1 - \Phi_i} \underline{X}_i + \sum_{\{i: D_i=1\}} (Y_i - \underline{X}_i^\top \underline{\beta}) \underline{X}_i = \underline{0}$$

which is equivalent to $(***)$. This shows that the $\underline{\beta}$ -limit of the EM algorithm is the MLE of $\underline{\beta}_0$.

(Exercise: Show that the solution to $\frac{\partial}{\partial \sigma^2} \frac{1}{n} \log \hat{L}_n(\underline{\theta}) = 0$ is the σ^2 -limit of the EM algorithm)

- Binary response models
- Multinomial response models
- Models with censoring/truncation

Suppose $[\underline{X}^\top Y]^\top$ is a random vector with

$$\begin{aligned} P[Y = 1|\underline{X}] &= g(\underline{X}) \\ P[Y = 0|\underline{X}] &= 1 - g(\underline{X}) \end{aligned}$$

where $g(\cdot)$ is a Borel measurable function into $(0,1)$. Different models of binary response involve different specifications of the link function $g(\cdot)$.

Typically one assumes a parametrization, i.e. $g(\underline{x})=g(\underline{x};\underline{\theta}_0)$ for some parameter $\underline{\theta}_0 \in \Theta$.

Typically use ML to estimate $\underline{\theta}_0$.

Where do binary response models come from?

- Standard answer: Involves a latent variable model.

Let Y^* be a continuous random variable that is only partially observed (i.e. Y^* is latent). In particular, assume that we observe

$$Y = \begin{cases} 1, & Y^* > 0 \\ 0, & Y^* \leq 0 \end{cases}$$

In addition, usually assume that Y^* is related to observable regressors \underline{X} via an index model.

$$Y^* = \underline{X}^\top \underline{\beta}_0 - U, \quad \dots (*)$$

where $\underline{\beta}_0$ is an unknown vector of constants and U is an unobservable random variable.

(We do not observe Y^* , only observe $Y \equiv 1\{Y^* > 0\}$.)

From (*) we have

$$\begin{aligned} P[Y = 1|\underline{X}] &= P[Y^* > 0] \\ &= P[\underline{X}^\top \underline{\beta}_0 - U > 0] \\ &= F_U(\underline{X}^\top \underline{\beta}_0) \end{aligned}$$

where F_U is the cdf of U . In the context of (*), we see that different models of binary response are characterized by different distributions of U .

e.g. Suppose $U \sim \text{logistic}$, i.e. U is continuous on \mathbb{R} with density

$$f_U(u) = \frac{e^u}{(1 + e^u)^2}$$

which implies that

$$P[Y = 1|\underline{X}] = F_U(\underline{X}^\top \underline{\beta}_0) = \frac{e^{\underline{X}^\top \underline{\beta}_0}}{1 + e^{\underline{X}^\top \underline{\beta}_0}}$$

i.e. a logit model.

e.g. Suppose $U \sim N(0,1)$. Then $P[Y=1|X] = \Phi(\underline{X}^T \underline{\beta}_0)$ i.e. a probit model.

Identification of $\underline{\beta}_0$ in a latent index model

$$Y^* = \underline{X}^T \underline{\beta}_0 - U$$

What we typically impose are restrictions on the variance and mean of U_i .

Observe that $Y=1\{Y^*>0\}=1\{\underline{X}^T \underline{\beta}_0 - U > 0\}$, but $1\{\underline{X}^T \underline{\beta}_0 - U > 0\}=1\{c\underline{X}^T \underline{\beta}_0 - cU > 0\}$ for any $c>0$. This means that observations generated by parameter $\underline{\beta}_0$ and latent error U are indistinguishable from observations generated by parameter $c\underline{\beta}_0$ and latent error cU for any $c>0$.

This problem is solved by fixing the variance of the latent error.

e.g. In a logit model, assume that the latent error has variance $\pi^2/3$, i.e. the variance of a logistic random variable. So, if $U \sim \text{logistic}$ and $\text{Var}[cU] = \pi^2/3$, then $c=1$.

In addition, note that \forall constant d

$$1\{\underline{X}^T \underline{\beta}_0 - U > 0\} = 1\{d + \underline{X}^T \underline{\beta}_0 - U - d > 0\}$$

which shows that it is impossible to distinguish between observations generated by parameter $\underline{\beta}_0$ and latent error U and observations generated by parameter $[d \quad \underline{\beta}_0^T]^T$ and latent error $U-d$.

This problem can be solved by requiring the latent error to have mean zero.

e.g. If $U \sim N(0,1)$ and $E[U-d]=0$ then $d=0$.

An important, special case of a latent index model is an additive random-utility model (ARUM). Suppose a consumer has to pick between two alternatives (0 and 1). Observe

$$Y = \begin{cases} 0, & \text{if alternative 0 has higher utility} \\ 1, & \text{if alternative 1 has higher utility} \end{cases}$$

An ARUM specifies the utilities of alternatives 0,1, as

$$\begin{aligned} U_0 &= V_0 + \varepsilon_0 \\ U_1 &= V_1 + \varepsilon_1 \end{aligned}$$

Call V_0, V_1 the deterministic components of utility; $\varepsilon_0, \varepsilon_1$ the random components of utility, then

$$\begin{aligned} P[Y = 1] &= P[U_1 \geq U_0] \\ &= P[V_1 - V_0 \geq \varepsilon_0 - \varepsilon_1] \\ &= F_{\varepsilon_0 - \varepsilon_1}(V_1 - V_0) \end{aligned}$$

where $F_{\varepsilon_0 - \varepsilon_1}$ is the cdf of $\varepsilon_0 - \varepsilon_1$. Frequently assume that $V_1 - V_0 = \underline{X}^T \underline{\beta}_0$ for some observable random vector \underline{X} and some unobservable constant vector $\underline{\beta}_0$. Identification of $\underline{\beta}_0$ is guaranteed by assuming a fixed value (e.g. 1 or $\pi^2/3$) for the difference in random utility components and also a mean of zero for the difference in random utility components.

In the context of binary choice modeling often encounter the assumption that $\varepsilon_0, \varepsilon_1$ have type-I extreme-value distributions, i.e. ε_j ($j=0,1$) is continuous on \mathbb{R} with density

$$f_\varepsilon(\varepsilon') = e^{-\varepsilon'} e^{-e^{-\varepsilon'}}$$

which implies a cdf $F_\varepsilon(\varepsilon') = e^{-e^{-\varepsilon'}}$.

Can show that if $\varepsilon_0, \varepsilon_1$ are independent and type-I extreme-value, then $\varepsilon_0 - \varepsilon_1 \sim \text{Logistic}$, and the binary-choice model $P[Y=1] = F_{\varepsilon_0 - \varepsilon_1}(V_1 - V_0)$ is a logit.

Often specifies ARUMs where the deterministic component of utility for alternative $j \in \{0, 1\}$ is

$$V_{ij} = \underline{Z}_{ij}^T \underline{\alpha}_j + \underline{W}_i^T \underline{\gamma}_j$$

for individual i .

In this case, \underline{Z}_{ij} are regressors that vary across alternatives (e.g. interactions between individual characteristics and product attributes) and \underline{W}_i are individual characteristics that are fixed across alternatives (e.g. income/gender/schooling). Then the binary-choice model has the form

$$P[Y_i = 1 | \underline{Z}_{i0}, \underline{Z}_{i1}, \underline{W}_i] = F_{\varepsilon_0 - \varepsilon_1}(\underline{Z}_{i1}^T \underline{\alpha}_1 - \underline{Z}_{i0}^T \underline{\alpha}_0 + \underline{W}_i^T (\underline{\gamma}_1 - \underline{\gamma}_0))$$

which show that $\underline{\gamma}_1, \underline{\gamma}_0$ cannot be separately identified.

General problem of MLE of a binary response model

Suppose we observe a random sample $[\underline{X}_i^T \ Y_i]^T$ ($i=1, \dots, n$) from a model with

$$\begin{aligned} P[Y_i = 1 | \underline{X}_i] &= G(\underline{X}_i^T \underline{\beta}_0) \\ P[Y_i = 0 | \underline{X}_i] &= 1 - G(\underline{X}_i^T \underline{\beta}_0) \end{aligned}$$

for some constant vector $\underline{\beta}_0 \in \mathbb{R}^d$ and where $G: \mathbb{R} \rightarrow (0, 1)$ is a positive, non-decreasing and continuously differentiable function on \mathbb{R} .

Then (obviously) Y_1, \dots, Y_n are discrete with probability mass function

$$f(y | \underline{X}_i; \underline{\beta}_0) = \begin{cases} G(\underline{X}_i^T \underline{\beta}_0)^y (1 - G(\underline{X}_i^T \underline{\beta}_0))^{1-y}, & y \in \{0, 1\} \\ 0, & \text{otherwise} \end{cases}$$

Evaluate the pmf at an actual observation Y_i ($i \in \{1, \dots, n\}$) replace $\underline{\beta}_0$ with an arbitrary $\underline{\beta} \in \mathbb{R}^d$ to get

$$\log f(Y_i | \underline{X}_i; \underline{\beta}) = Y_i \log G(\underline{X}_i^T \underline{\beta}) + (1 - Y_i) \log (1 - G(\underline{X}_i^T \underline{\beta}))$$

which implies the log likelihood

$$\begin{aligned} \hat{L}_n(\underline{\beta}) &= \sum_{i=1}^n \log f(Y_i | \underline{X}_i; \underline{\beta}) \\ &= \sum_{i=1}^n [Y_i \log G(\underline{X}_i^T \underline{\beta}) + (1 - Y_i) \log (1 - G(\underline{X}_i^T \underline{\beta}))] \end{aligned}$$

By differentiating

$$\frac{1}{n} \frac{\partial}{\partial \underline{\beta}^T} \hat{L}_n(\underline{\beta}) = \frac{1}{n} \sum_{i=1}^n \left(\frac{Y_i}{G(\underline{X}_i^T \underline{\beta})} G'(\underline{X}_i^T \underline{\beta}) \underline{X}_i - \frac{1 - Y_i}{1 - G(\underline{X}_i^T \underline{\beta})} G'(\underline{X}_i^T \underline{\beta}) \underline{X}_i \right)$$

The MLE $\hat{\beta}_n$ is given by

$$\frac{1}{n} \frac{\partial}{\partial \underline{\beta}^\top} \hat{L}_n(\underline{\beta}) \Big|_{\underline{\beta}=\hat{\beta}_n} = \underline{0} \Leftrightarrow \sum_{i=1}^n \frac{Y_i - G(\underline{X}_i^\top \hat{\beta}_n)}{G(\underline{X}_i^\top \hat{\beta}_n) (1 - G(\underline{X}_i^\top \hat{\beta}_n))} G'(\underline{X}_i^\top \hat{\beta}_n) \underline{X}_i = 0, \quad \dots (**)$$

We can see that in general, there is no closed form for $\hat{\beta}_n$.

Note: If G is the logistic or $N(0,1)$ -cdf then $\hat{L}_n(\underline{\beta})$ is globally (exercise) and Newton–Raphson method converges quickly.

When is the MLE $\hat{\beta}_n$ consistent?

Need the “success probability” $G(\underline{X}_i^\top \underline{\beta}_0)$ to be correctly specified.

Since Y_i is Bernoulli, we have

$$E[Y_i | \underline{X}_i] = P[Y_i = 1 | \underline{X}_i]$$

so we need $P[Y_i = 1 | \underline{X}_i] = G(\underline{X}_i^\top \underline{\beta}_0)$ if

$$\frac{1}{n} E \left[\frac{\partial}{\partial \underline{\beta}^\top} \hat{L}_n(\underline{\beta}) \Big|_{\underline{\beta}=\underline{\beta}_0} \right] = \underline{0} \Leftrightarrow E \left[\frac{E[Y_1 | \underline{X}_1] - G(\underline{X}_1^\top \underline{\beta}_0)}{G(\underline{X}_1^\top \underline{\beta}_0) (1 - G(\underline{X}_1^\top \underline{\beta}_0))} G'(\underline{X}_1^\top \underline{\beta}_0) \underline{X}_1 \right] = \underline{0}$$

Observe a random sample

$$[\underline{X}_i^T \ Y_i]^T, \quad i = 1, \dots, n$$

from a model with

$$\begin{aligned} P[Y_i = 1 | \underline{X}_i] &= G(\underline{X}_i^T \underline{\beta}_0) \\ P[Y_i = 0 | \underline{X}_i] &= 1 - G(\underline{X}_i^T \underline{\beta}_0) \end{aligned}$$

for some constant vector $\underline{\beta}_0 \in \mathbb{R}^d$ and some non-decreasing continuously differentiable function $G: \mathbb{R} \rightarrow (0,1)$. Then the pmf of each Y_i is

$$f(y | \underline{X}_i; \underline{\beta}_0) = \begin{cases} G(\underline{X}_i^T \underline{\beta}_0)^y (1 - G(\underline{X}_i^T \underline{\beta}_0))^{1-y}, & y \in \{0,1\} \\ 0, & \text{otherwise} \end{cases}$$

which implies the log-likelihood

$$\begin{aligned} \hat{L}_n(\underline{\beta}) &\equiv \sum_{i=1}^n \log f(Y_i | \underline{X}_i; \underline{\beta}) \\ &= \sum_{i=1}^n \left[Y_i \log G(\underline{X}_i^T \underline{\beta}) + (1 - Y_i) \log (1 - G(\underline{X}_i^T \underline{\beta})) \right] \\ \Rightarrow \frac{1}{n} \frac{\partial}{\partial \underline{\beta}^T} \hat{L}_n(\underline{\beta}) &= \frac{1}{n} \sum_{i=1}^n \left(\frac{Y_i}{G(\underline{X}_i^T \underline{\beta})} - \frac{1 - Y_i}{1 - G(\underline{X}_i^T \underline{\beta})} \right) G'(\underline{X}_i^T \underline{\beta}) \underline{X}_i \end{aligned}$$

The MLE $\hat{\underline{\beta}}_n$ is implicitly defined as

$$\left. \frac{1}{n} \frac{\partial}{\partial \underline{\beta}^T} \hat{L}_n(\underline{\beta}) \right|_{\underline{\beta} = \hat{\underline{\beta}}_n} = \underline{0} \Leftrightarrow \sum_{i=1}^n \frac{Y_i - G(\underline{X}_i^T \hat{\underline{\beta}}_n)}{G(\underline{X}_i^T \hat{\underline{\beta}}_n) (1 - G(\underline{X}_i^T \hat{\underline{\beta}}_n))} G'(\underline{X}_i^T \hat{\underline{\beta}}_n) \underline{X}_i = \underline{0}$$

($\hat{\underline{\beta}}_n$ in general has no closed form)

Note: With logit or probit, $G(\cdot)$ is globally concave and Newton–Raphson will converge quickly.

Consistency requires that $G(\cdot)$ be correctly specified: This has to do with the necessary condition that

$$E \left[\left. \frac{\partial}{\partial \underline{\beta}^T} \hat{L}_n(\underline{\beta}) \right|_{\underline{\beta} = \underline{\beta}_0} \right] = \underline{0}$$

which implies that

$$E[Y_i | \underline{X}_i] = P[Y_i = 1 | \underline{X}_i] = G(\underline{X}_i^T \underline{\beta}_0)$$

Asymptotic Normality in the standard case requires that $G(\cdot)$ to be correctly specified and twice continuously differentiable.

Under these assumptions, we can make the standard argument that the MLE is \sqrt{n} -consistent and asymptotically Normal with

$$\sqrt{n}(\hat{\beta}_n - \beta_0) \xrightarrow{d} N_d \left(\underline{0}, \left(\text{plim} \left(-\frac{1}{n} \frac{\partial^2}{\partial \underline{\beta} \partial \underline{\beta}^\top} \hat{L}_n(\underline{\beta}) \Big|_{\underline{\beta}=\underline{\beta}_0} \right) \right)^{-1} \right)$$

Note that

$$\begin{aligned} \frac{1}{n} \frac{\partial^2}{\partial \underline{\beta} \partial \underline{\beta}^\top} \hat{L}_n(\underline{\beta}) &= \frac{1}{n} \sum_{i=1}^n \left(-\frac{Y_i}{G^2(\underline{X}_i^\top \underline{\beta})} G'(\underline{X}_i^\top \underline{\beta})^2 \underline{X}_i \underline{X}_i^\top + \frac{Y_i}{G(\underline{X}_i^\top \underline{\beta})} G''(\underline{X}_i^\top \underline{\beta}) \underline{X}_i \underline{X}_i^\top \right. \\ &\quad \left. + \frac{1-Y_i}{(1-G(\underline{X}_i^\top \underline{\beta}))^2} G'(\underline{X}_i^\top \underline{\beta})^2 \underline{X}_i \underline{X}_i^\top - \frac{1-Y_i}{1-G(\underline{X}_i^\top \underline{\beta})} G''(\underline{X}_i^\top \underline{\beta}) \underline{X}_i \underline{X}_i^\top \right) \end{aligned}$$

This implies that

$$\begin{aligned} \frac{1}{n} \frac{\partial^2}{\partial \underline{\beta} \partial \underline{\beta}^\top} \hat{L}_n(\underline{\beta}) \Big|_{\underline{\beta}=\underline{\beta}_0} &\xrightarrow{p} -E \left[\frac{G'(\underline{X}_1^\top \underline{\beta}_0)^2}{G(\underline{X}_1^\top \underline{\beta}_0)} \underline{X}_1 \underline{X}_1^\top + G''(\underline{X}_1^\top \underline{\beta}_0) \underline{X}_1 \underline{X}_1^\top \right. \\ &\quad \left. + \frac{G'(\underline{X}_1^\top \underline{\beta}_0)^2}{1-G(\underline{X}_1^\top \underline{\beta}_0)} \underline{X}_1 \underline{X}_1^\top - G''(\underline{X}_1^\top \underline{\beta}_0) \underline{X}_1 \underline{X}_1^\top \right] \\ &= -E \left[\frac{G'(\underline{X}_1^\top \underline{\beta}_0)^2}{G(\underline{X}_1^\top \underline{\beta}_0) (1-G(\underline{X}_1^\top \underline{\beta}_0))} \underline{X}_1 \underline{X}_1^\top \right] \end{aligned}$$

Therefore a consistent estimator of the asymptotic variance of $\hat{\beta}_n$ is

$$\hat{V}[\hat{\beta}_n] = \left[\frac{1}{n} \sum_{i=1}^n \frac{G'(\underline{X}_i^\top \hat{\beta}_n)^2}{G(\underline{X}_i^\top \hat{\beta}_n) (1-G(\underline{X}_i^\top \hat{\beta}_n))} \underline{X}_i \underline{X}_i^\top \right]^{-1}$$

More on logit: Logit satisfies

$$G(\underline{X}^\top \underline{\beta}) = \frac{e^{\underline{X}^\top \underline{\beta}}}{1 + e^{\underline{X}^\top \underline{\beta}}}$$

The corresponding FOC are

$$\frac{1}{n} \sum_{i=1}^n \left(Y_i - \frac{e^{\underline{X}_i^\top \underline{\beta}}}{1 + e^{\underline{X}_i^\top \underline{\beta}}} \right) \underline{X}_i = \underline{0}$$

because $G'(\underline{X}^\top \underline{\beta}) = G(\underline{X}^\top \underline{\beta}) (1 - G(\underline{X}^\top \underline{\beta}))$

If \underline{X} includes a constant term (i.e. $\underline{X}^T \underline{\beta}$ includes an intercept), then the FOCs imply that

$$\sum_{i=1}^n \left(Y_i - \frac{e^{\underline{X}_i^T \underline{\beta}}}{1 + e^{\underline{X}_i^T \underline{\beta}}} \right) = 0 \Leftrightarrow \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i = \frac{1}{n} \sum_{i=1}^n G(\underline{X}_i^T \underline{\beta}_n)$$

i.e. the sample proportion \bar{Y} is equal to the in-sample predicted probability.

Note also that

$$G(\underline{X}^T \underline{\beta}) = \frac{e^{\underline{X}^T \underline{\beta}}}{1 + e^{\underline{X}^T \underline{\beta}}} \Leftrightarrow e^{\underline{X}^T \underline{\beta}} = \frac{G(\underline{X}^T \underline{\beta})}{1 - G(\underline{X}^T \underline{\beta})} \Leftrightarrow \log \frac{G(\underline{X}^T \underline{\beta})}{1 - G(\underline{X}^T \underline{\beta})} = \log \frac{P[Y = 1|\underline{X}]}{P[Y = 0|\underline{X}]} = \underline{X}^T \underline{\beta}$$

Therefore if β_j is the j th component of $\underline{\beta}$ and X_j is the j th component of \underline{X} then β_j is the change in the log-odds ratio. When X_j increases by one unit and all other components of \underline{X} are constant

e.g. Suppose $\Delta X_j = 1$. Then $e^{\underline{X}^T \underline{\beta}}$ increases to $e^{\underline{X}^T \underline{\beta} + \beta_j} = e^{\underline{X}^T \underline{\beta}} e^{\beta_j}$, i.e. the odds ratio increases by e^{β_j} . So if $\beta_j = .10$, then $\Delta X_j = 1$ multiplies the odds ratio by $e^{.10} \approx 1.105$, i.e. a 10.5% increase.

More on the Probit model

Probit specifies $G(\underline{X}^T \underline{\beta}) = \Phi(\underline{X}^T \underline{\beta})$ where $\Phi(\cdot)$ is the $N(0,1)$ -cdf. The corresponding FOCs are

$$\frac{1}{n} \sum_{i=1}^n \frac{\varphi(\underline{X}_i^T \underline{\beta})}{\Phi(\underline{X}_i^T \underline{\beta}) (1 - \Phi(\underline{X}_i^T \underline{\beta}))} (Y_i - \Phi(\underline{X}_i^T \underline{\beta})) \underline{X}_i = \underline{0}$$

(compare with the logit, OLS) Also,

$$\frac{\partial}{\partial X_j} G(\underline{X}^T \underline{\beta}) = \varphi(\underline{X}^T \underline{\beta}) \beta_j$$

So β_j is not intuitively interpretable.

What about OLS? OLS in this context assumes a linear probability model, i.e. the model

$$\begin{aligned} P[Y = 1|\underline{X}] &= \underline{X}^T \underline{\beta}_0 \\ P[Y = 0|\underline{X}] &= 1 - \underline{X}^T \underline{\beta}_0 \end{aligned}$$

which implies $E[Y|\underline{X}] = \underline{X}^T \underline{\beta}_0$.

OLS will work as expected, but generally yield predicted probabilities $\underline{X}^T \hat{\underline{\beta}}_{OLS}$ that are not in the interval $[0,1]$.

If OLS is to be used, then heteroskedasticity-consistent standard errors should be used, since

$$\text{Var}[Y|\underline{X}] = \underline{X}^T \underline{\beta}_0 (1 - \underline{X}^T \underline{\beta}_0)$$

which varies with \underline{X} .

How to choose between logit or probit?

In the first instance, we should choose depending on how confident we are in the assumed functional form of $G(\cdot)$, i.e.

$$\text{if } P[Y = 1|\underline{X}] = F(\underline{X}^T \underline{\beta}_0)$$

where $F(\cdot)$ is the logistic or $N(0,1)$ -cdf, then logit or probit should be used as appropriately.

In the second instance, misspecification of $G(\cdot)$ may not be serious.

Ruud (1983, EMA): Sufficient condition for scale consistency in this context is that for each j

$$E[X_j | \underline{X}^T \underline{\beta}_0 = z] \text{ is linear in } z$$

Scale consistency: i.e. $\hat{\beta}_n \xrightarrow{p} c\beta_0$ for some constant c

Unfortunately, the condition that $E[X_j | \underline{X}^T \underline{\beta}_0 = z]$ is linear in z is hard to verify.

In practice: Usually little difference between predicted probabilities coming from logit or probit, with the exception of predicted probabilities near 0 or 1.

Diagnostic measures

a) Pseudo- R^2

A popular measure is the McFadden (1974, *Frontiers in Econometrics*) Pseudo- R^2 , i.e.

$$R_M^2 \equiv 1 - \frac{\hat{L}_n(\hat{\beta}_n)}{\hat{L}_n(\bar{Y})} = 1 - \frac{\sum_{i=1}^n [Y_i \log G(\underline{X}_i^T \hat{\beta}_n) + (1 - Y_i) \log (1 - G(\underline{X}_i^T \hat{\beta}_n))]}{n\bar{Y} \log \bar{Y} + (1 - \bar{Y}) \log(1 - \bar{Y})}$$

Suggestion: Compare logit/probit/OLS with a “distribution-free” method that allows for misspecification of $G(\cdot)$. Refer Klein and Spady (1993, EMA), “Semiparametric MLE” and Ichimura (1993, JEM), “Semiparametric Least Squares” for “distribution-free” cases.

Values of R_M^2 close to one indicate a fitted binary response model that predicts $Y=1$ with probability close to $G(\underline{X}_i^T \hat{\beta}_n)$ and $Y=0$ with probability close to $1 - G(\underline{X}_i^T \hat{\beta}_n)$.

b) Quality of predicted outcomes

Let $\hat{Y} \in \{0,1\}$ be a predicted value of Y . Can plot a receiver operating characteristics (ROC) curve by setting $\hat{Y}=1$ if $G(\underline{X}_i^T \hat{\beta}_n) > \tau$ for some $\tau \in (0,1)$; $\hat{Y}=0$ otherwise.

The ROC curve plots

$$\frac{\#\{Y = 1 \text{ values correctly predicted}\}}{\#\{Y = 0 \text{ values correctly predicted}\}} \text{ as a function of } \tau.$$

So if $\tau=1$, then all $Y=0$ values are correctly predicted, and all $Y=1$ values are incorrectly predicted, i.e. the ROC curve is equal to 0 at $\tau=1$; similarly, the ROC curve is ∞ at $\tau=0$. Ideally want a ROC curve close to 1 for all $\tau \in (0,1)$.

Example

- Sample of $n=630$ recreational fishers
- For each $i \in \{1, \dots, 630\}$

$$\text{Let } Y_i = \begin{cases} 1, & \text{if fishing from a charter boat} \\ 0, & \text{if fishing from a pier} \end{cases}, \quad X_i \equiv \log \text{RELP}_i \equiv \log \left(\frac{P_{\text{Charter},i}}{P_{\text{Pier},i}} \right)$$

- Obviously: Expect the incidence of charter boat fishing to be decreasing in X_i
- We observe $\bar{Y} \equiv \frac{1}{n} \sum Y_i = .717$
- For observations i with $Y_i=1$, the average RELP was $75/121 \approx .6198$. So price apparently has the expected effect
- We also find that for observations i with $Y_i=0$, the average RELP was $110/31 \approx 3.5484$

How to model the relationship between Y_i and X_i ?

- OLS. Regress Y_i on a constant and X_i (cannot constrain predicted choice probabilities to be between 0 and 1): Intercept=.784, Slope=.243
- Logit (the assumption is that

$$P[Y_i = 1|X_i] = \frac{e^{\beta_{01} + \beta_{02}X_i}}{1 + e^{\beta_{01} + \beta_{02}X_i}}$$

for some constants β_{01} and β_{02} : $\hat{\beta}_{01} = 2.035, \hat{\beta}_{02} = -1.823$.

The implied estimate of the marginal effect of X_i on the choice of probability is

$$\frac{\hat{\beta}_{02} e^{\hat{\beta}_{01} + \hat{\beta}_{02}X_i}}{(1 + e^{\hat{\beta}_{01} + \hat{\beta}_{02}X_i})^2} < 0$$

for all i as expected. The estimated average marginal effect

$$\bar{Y}(1 - \bar{Y})\hat{\beta}_{02} = -.370$$

- Probit (the assumption is that

$$P[Y_i = 1|X_i] = \Phi(\beta_{01} + \beta_{02}X_i)$$

for constants β_{01}, β_{02} : $\hat{\beta}_{01} = 1.194, \hat{\beta}_{02} = -1.056$. The estimated marginal effect is

$$\Phi(\hat{\beta}_{01} + \hat{\beta}_{02}X_i)\hat{\beta}_{02} < 0$$

for all i as expected.

In this example, the logit and probit predicted choice probabilities are similar for all X_i , which is often.

Multinomial response models

Setup: m mutually exclusive alternatives. The dependent variable for each observation i is conditionally multinomial given regressors \underline{X}_i , i.e.

$$Y_i = j, \quad \text{if } i \text{ chooses alternative } j \in \{1, \dots, m\}$$

e.g. Mode of commuting

$$Y_i = \begin{cases} 1, & \text{if drives own car} \\ 2, & \text{if bus} \\ 3, & \text{if walk} \\ 4, & \text{if bike} \end{cases}$$

Now define binary indicators

$$Y_{ij} = 1\{Y_i = j\}, \quad \text{for } j \in \{1, \dots, m\}$$

Therefore, the (multinomial) density for a single observation is

$$f_Y(y) \equiv \prod_{j=1}^m P_{ij}^{y_{ij}}$$

where $P_{ij} = P[Y_i=j|\underline{X}_i] = P[Y_{ij}=1|\underline{X}_i]$. Typically model each P_{ij} as

$$P_{ij} = G_j(\underline{X}_i, \underline{\beta}_0)$$

for some unknown constant $\underline{\beta}_0$. Call $G_j(\underline{x}, \underline{\beta})$ ($j=1, \dots, m$) the multinomial link functions. Different specifications of link function correspond to different models of multinomial response.

Maximum likelihood estimation

Suppose for individual i and alternative j

$$P_{ij} = P[Y_{ij} = 1|\underline{X}_i] = G_j(\underline{X}_i, \underline{\beta}_0)$$

Then the log-likelihood is

$$\hat{L}_n(\underline{\beta}) = \sum_{i=1}^n \sum_{j=1}^m Y_{ij} \log G_j(\underline{X}_i, \underline{\beta})$$

(assuming a sample of n individuals)

In many applications, individual i faces a mean of m_i alternatives, so

$$\hat{L}_n(\underline{\beta}) = \sum_{i=1}^n \sum_{j=1}^{m_i} Y_{ij} \log G_j(\underline{X}_i, \underline{\beta})$$

FOCs given by

$$\left. \frac{1}{n} \frac{\partial}{\partial \underline{\beta}^\top} \hat{L}_n(\underline{\beta}) \right|_{\underline{\beta}=\hat{\underline{\beta}}_n} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{m_i} \frac{Y_{ij}}{G_j(\underline{X}_i, \hat{\underline{\beta}}_n)} \frac{\partial}{\partial \underline{\beta}^\top} G_j(\underline{X}_i, \underline{\beta}) \Big|_{\underline{\beta}=\hat{\underline{\beta}}_n} = 0$$

Sufficient for consistency is the score identity, which holds if $P_{ij}=G_j(\underline{X}_i, \underline{\beta}_0)$, i.e. correct specification of each link function (Exercise: Under correct specification $\sum_{j=1}^{m_i} G_j(\underline{X}_i, \underline{\beta}_0) = 1$ for each i).

Note: The standard consistency argument for MLE on the score identity relies on $G_j(\underline{X}_i, \underline{\beta}_0)$ being correctly specified and $G_j(\underline{X}_i, \underline{\beta})$ being continuously differentiable at $\underline{\beta} = \underline{\beta}_0$.

If link functions are correctly specified and twice continuously differentiable at $\underline{\beta} = \underline{\beta}_0$, then the standard argument shows that the MLE $\hat{\underline{\beta}}_n$ is \sqrt{n} -consistent and asymptotically Normal with

$$\sqrt{n}(\hat{\underline{\beta}}_n - \underline{\beta}_0) \xrightarrow{d} N\left(0, \text{plim} \left[-\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{m_i} \left(\frac{1}{G_j(\underline{X}_i, \underline{\beta}_0)} \frac{\partial G_j(\underline{X}_i, \underline{\beta})}{\partial \underline{\beta}^\top} \bigg|_{\underline{\beta}=\underline{\beta}_0} \frac{\partial G_j(\underline{X}_i, \underline{\beta})}{\partial \underline{\beta}} \bigg|_{\underline{\beta}=\underline{\beta}_0} - \frac{\partial^2 G_j(\underline{X}_i, \underline{\beta})}{\partial \underline{\beta} \partial \underline{\beta}^\top} \bigg|_{\underline{\beta}=\underline{\beta}_0} \right) \right]^{-1} \right)$$

Exercise: Fill in the details.

Next: MLE is convenient for iid observations. When this assumption is implausible (say independence without identical marginals), can adopt moment-based estimator of the form

$$\sum_{i=1}^n \sum_{j=1}^{m_i} (Y_{ij} - G_j(\underline{X}_i, \hat{\underline{\beta}}_n)) \underline{w}_i = \underline{0}, \quad (*)$$

where \underline{w}_i has the same dimension as $\underline{\beta}_0$, but does not depend on Y_{ij} .

Under further (but standard) conditions, $\hat{\underline{\beta}}_n$ is \sqrt{n} -consistent and asymptotically Normal with an asymptotic variance that depends on \underline{W}_i (Exercise: Derive the asymptotic distribution of $\sqrt{n}(\hat{\underline{\beta}}_n - \underline{\beta}_0)$ under these conditions). GMM theory can help select the efficient setting of \underline{W}_i .

We've seen that different specifications of link function imply different models of multinomial response.

Can also get different multinomial response models for a given link function by allowing regressors or parameters to vary across alternatives.

1) Alternative-varying regressors

For individual i , suppose the regressors are

$$\underline{X}_i = [\underline{X}_{i1}^\top \quad \cdots \quad \underline{X}_{im}^\top]^\top$$

The link function for alternative j has a multi-index form given by

$$G_j(\underline{X}_i, \underline{\beta}_0) = G_j(\underline{X}_{i1}^\top \underline{\beta}_{01}, \dots, \underline{X}_{im}^\top \underline{\beta}_{0m})$$

assuming $\underline{\beta}_0$ is invariant across alternatives.

2) Alternative-invariant regressors

This is typically coupled with alternative-varying parameters. Then the link function also has a multi-index form

$$G_j(\underline{X}_i, \underline{\beta}_0) = G_j(\underline{X}_i^\top \underline{\beta}_{01}, \dots, \underline{X}_i^\top \underline{\beta}_{0m})$$

where $\underline{\beta}_0 = [\underline{\beta}_{01}^\top \quad \cdots \quad \underline{\beta}_{0m}^\top]^\top$

In this case, typically normalize

$$\underline{\beta}_{01} = \underline{0}$$

(Important: Software for multinomial response models requires one to specify alternative-varying regressors with alternative-invariant parameters or vice versa)

Recommendation: Use software for alternative varying regressors. Then if alternative-invariant regressors seem more plausible, then write

$$X_{ij} \equiv \left[\underbrace{0^T \cdots 0^T}_{j-1 \times d} \quad \underbrace{X_i^T}_{1 \times d} \quad \underbrace{0^T \cdots 0^T}_{n-j} \right]^T$$

Where if X_i is a d-variate, $\underline{0}$ is a $(d \times 1)$ -vector of zeroes.

In addition, need to set

$$\underline{\beta} \equiv [\underline{0}^T \quad \underline{\beta}_2^T \quad \cdots \quad \underline{\beta}_m^T]^T$$

where each $\underline{\beta}_j$ ($j=1, \dots, m$) is $d \times 1$, and $\underline{\beta}_1 = \underline{0}$ is a normalization.

The simplest category of multinomial response model involves logistic link function (Luce, 1959).

Several different variants

- a) Alternative-varying regressors, i.e. conditional logit (CL), we have

$$G_j(\underline{X}_i, \underline{\beta}_0) = \frac{e^{\underline{X}_{ij}^T \underline{\beta}_0}}{1 + e^{\underline{X}_{ij}^T \underline{\beta}_0}}, \quad j = 1, \dots, m$$

- b) Alternative-invariant regressors, i.e. multi-nomial logit (MNL)

$$G_j(\underline{X}_i, \underline{\beta}_0) = \frac{e^{\underline{X}_i^T \underline{\beta}_{0j}}}{1 + e^{\underline{X}_i^T \underline{\beta}_{0j}}}, \quad \text{for } j \in \{1, \dots, m\} \text{ and where } \underline{\beta}_{01} = \underline{0}$$

- c) A combination of alternative-varying and alternative-invariant regressors, i.e. mixed logit: Let $\underline{X}_i \equiv [\underline{X}_{i1}^T \quad \cdots \quad \underline{X}_{im}^T]^T$ be the vector of alternative-varying regressors; let \underline{W}_i be the vector of alternative-invariant regressors. Let $\underline{\gamma}_0 \equiv [\underline{\gamma}_{01}^T \quad \cdots \quad \underline{\gamma}_{0m}^T]^T$, then

$$P_{ij} = G_j(\underline{X}_i, \underline{W}_i, \underline{\beta}_0, \underline{\gamma}_0) = \frac{e^{\underline{X}_{ij}^T \underline{\beta}_0 + \underline{W}_i^T \underline{\gamma}_{0j}}}{\sum_{k=1}^m e^{\underline{X}_{ik}^T \underline{\beta}_0 + \underline{W}_i^T \underline{\gamma}_{0k}}}$$

Luce (1959): Conditional Logit (CL)

$$P[Y_i = j | X_i] = G_j(X_i; \beta_0) \equiv \frac{e^{X_{ij}^T \beta_0}}{\sum_{k=1}^m e^{X_{ik}^T \beta_0}}, \quad \text{individual } i, \quad j \in \{1, \dots, m\}$$

Multinomial Logit (MNL)

$$P[Y_i = j | X_i] = G_j(X_i; \beta_0) \equiv \frac{e^{X_i^T \beta_{0j}}}{\sum_{k=1}^m e^{X_i^T \beta_{0k}}}, \quad \text{individual } i, \quad j \in \{1, \dots, m\}$$

MNL requires normalization, usually $\beta_{01} = 0$.

Maximum-likelihood Estimation of CL and MNL; for CL

$$\frac{\partial}{\partial \beta^T} G_j(X_i, \beta) = G_j(X_i, \beta) (X_{ij} - \bar{X}_i), \quad \text{where } \bar{X}_i = \sum_{j=1}^m G_j(X_i, \beta) X_{ij}$$

i.e. conditional probability weighted average of regressors. This implies that the FOCs are

$$\sum_{i=1}^n \sum_{j=1}^m Y_{ij} (X_{ij} - \bar{X}_i) = \underline{0}, \quad (\text{Exercise})$$

In addition, the MLE of a CL satisfies

$$\sqrt{n}(\hat{\beta}_{CL} - \beta_0) \xrightarrow{d} N \left(\underline{0}, \text{plim} \left(\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m G_j(X_i, \beta_0) (X_{ij} - \bar{X}_i) (X_{ij} - \bar{X}_i)^T \right)^{-1} \right)$$

(Exercise: Fill in details) For MNL, we have for each $j \in \{1, \dots, m\}$

$$\frac{\partial}{\partial \beta_k^T} G_j(X_i, \beta) = G_j(X_i, \beta) (\delta_{ijk} - G_k(X_i, \beta)) X_i$$

Where $\delta_{ijk} = 1 \{j=k\}$. So the FOCs have the form

$$\sum_{i=1}^n (Y_{ij} - G_j(X_i, \beta)) X_i = \underline{0}, \quad \text{for each } j \in \{1, \dots, m\}$$

Also, the MNL MLE satisfies

$$\sqrt{n}(\hat{\beta}_{\text{MNL}} - \beta_0) \xrightarrow{d} N\left(0, \text{plim}\left(-\frac{1}{n} \frac{\partial^2 \hat{L}_n(\beta)}{\partial \beta \partial \beta^\top} \Big|_{\beta=\beta_0}\right)^{-1}\right)$$

where the (j,k)-block of $\text{plim}\left(-\frac{1}{n} \frac{\partial^2 \hat{L}_n(\beta)}{\partial \beta \partial \beta^\top} \Big|_{\beta=\beta_0}\right)$ is

$$\text{plim} \frac{1}{n} \sum_{i=1}^n G_j(\underline{X}_i, \beta_0) (\delta_{ijk} - G_k(\underline{X}_i, \beta_0)) \underline{X}_i \underline{X}_i^\top, \quad \text{where } j, k \in \{1, \dots, m\}$$

(Exercise: Fill in details)

Next: Consider the interpretation of CL and MNL model parameters

CL model: What is the effect on $G_j(\underline{X}_i, \beta_0)$ of an infinitesimal change in \underline{X}_{ik} ?

$$\frac{\partial G_j(\underline{X}_i, \beta_0)}{\partial \underline{X}_{ik}^\top} = G_j(\underline{X}_i, \beta_0) (\delta_{ijk} - G_k(\underline{X}_i, \beta_0)) \underline{\beta}_0$$

So, if the e-th component of $\underline{\beta}_0$ is positive, an infinitesimal change in the e-th component of \underline{X}_{ik} increases the probability of alternative k and decreases the probability of every other alternative.

MNL model: What is the effect on $G_j(\underline{X}_i, \beta_0)$ of an infinitesimal change in \underline{X}_i ?

The answer is summarized by

$$\frac{\partial G_j(\underline{X}_i, \beta_0)}{\partial \underline{X}_i^\top} = G_j(\underline{X}_i, \beta_0) (\underline{\beta}_{0j} - \bar{\beta}_{0i}), \quad \text{where } \bar{\beta}_{0i} = \sum_{j=1}^m G_j(\underline{X}_i, \beta_0) \underline{\beta}_{0j}$$

(Exercise: Verify this)

Note: The sign of this effect is ambiguous.

An analytically convenient feature of CL or MNL is the possibility of rewriting either model in terms of a logit.

Take MNL: MNL is equivalent to a series of pairwise comparisons of all but one alternative to a “baseline category,” i.e. an alternative whose coefficient vector, say $\underline{\beta}_{01}$, is zero. To see this, note that

$$P[Y_i = j | Y_i \in \{j, k\}, \underline{X}_i] = \frac{G_j(\underline{X}_i, \beta_0)}{G_j(\underline{X}_i, \beta_0) + G_k(\underline{X}_i, \beta_0)} = \frac{e^{\underline{X}_i^\top \underline{\beta}_{0j}}}{e^{\underline{X}_i^\top \underline{\beta}_{0j}} + e^{\underline{X}_i^\top \underline{\beta}_{0k}}} = \frac{e^{\underline{X}_i^\top (\underline{\beta}_{0j} - \underline{\beta}_{0k})}}{1 + e^{\underline{X}_i^\top (\underline{\beta}_{0j} - \underline{\beta}_{0k})}}$$

i.e. a logit model with coefficient vector $\underline{\beta}_{0j} - \underline{\beta}_{0k}$.

Usually, set up alternative 1 as the baseline, so $\beta_{01} = 0$. In this case, for $j \in \{2, \dots, m\}$

$$P[Y_i = j | Y_i \in \{1, j\}, \underline{X}_i] = \frac{e^{\underline{X}_i^T \beta_{0j}}}{1 + e^{\underline{X}_i^T \beta_{0j}}}$$

i.e. a logit with coefficient vector β_{0j} . Thus, the relative risk of picking alternative j over alternative 1 is

$$\frac{P[Y_i = j | Y_i \in \{1, j\}, \underline{X}_i]}{P[Y_i = 1 | Y_i \in \{1, j\}, \underline{X}_i]} = e^{\underline{X}_i^T \beta_{0j}}$$

which implies that $e^{\beta_{0j\ell}}$ (where $\beta_{0j\ell}$ is the ℓ -th component of β_{0j}) is the instantaneous change in the relative risk when the ℓ -th component of \underline{X}_i of \underline{X}_i increases by one unit.

It follows that MNL coefficients only have a natural interpretation if there exists a logical base alternative.

e.g. Let $Y_i = j$ according to which of m modes of commuting individual I uses to get to work.

Depending on the audience/geography might make sense to see $Y_i = 1$ if individual I chooses to drive his/her own car.

Next: Consider CL. Can show that

$$P[Y_i = j | Y_i \in \{j, k\}, \underline{X}_i] = \frac{e^{(\underline{X}_{ij} - \underline{X}_{ik})^T \beta_0}}{1 + e^{(\underline{X}_{ij} - \underline{X}_{ik})^T \beta_0}}$$

i.e. a logit with regressors $\underline{X}_{ij} - \underline{X}_{ik}$. Therefore, should normalize w.r.t. the regressor values (say \underline{X}_{i1}) for the base alternative.

Important: The reducibility of CL or MNL to a series of $m-1$ logits is a serious limitation because it is implausible that individual choice over m alternatives can be reduced to a series of $m-1$ (pairwise comparisons).

e.g. “Red-bus, blue-bus” problem: Suppose $Y_i = 1, 2, 3$ as commuter I chooses to drive his own car to work ($Y_i = 1$), ride the red bus ($Y_i = 2$) or ride the blue bus ($Y_i = 3$).

Assume that alternatives 2 and 3 do not differ apart from the color of the bus, which is assumed to be irrelevant.

MNL or CL imply that conditional on \underline{X}_i , the probability of observing $Y_i = 3$ (red bus) is the same regardless of whether the choice set is $\{\text{car, red bus}\}$ ($m=2$) or $\{\text{car, red bus, blue bus}\}$ ($m=3$). This is unrealistic.

- Introducing an equivalent blue bus when a red bus is available should have the use of the red bus and have no effect on the popularity of driving

Essentially: CL and MNL both imply independence of irrelevant alternatives (IIA)

- The implausibility of IIA has led to an emphasis on multinomial response models that don't have this feature

e.g. nested logit, random parameters logit (RPL), multinomial probit (MNP)

- IIA is testable via the Hausman test: Hausman and McFadden (1984, EMA)

Consider the fundamental question of how to derive a multinomial response model for consumer choice: Easiest to start with an additive random utility (ARUM) model, where alternative $j \in \{1, \dots, m\}$ has utility

$$u_j = v_j + \varepsilon_j$$

where v_j is deterministic and ε_j is random. The alternative chosen is the one with the highest u_j . Suppose that

$$u_{ij} = v_{ij} + \varepsilon_{ij}, \quad \text{for individual } i, \quad j \in \{1, \dots, m\}$$

Then,

$$\begin{aligned} P[i \text{ chooses } j] &= P[Y_i = j] \\ &= P[u_{ij} + u_{ik}], \quad \forall j \neq k \\ &= P[\varepsilon_{ik} - \varepsilon_{ij} \leq v_{ij} - v_{ik}], \quad \forall j \neq k \\ &= F_{\varepsilon_{ik} - \varepsilon_{ij}}(v_{ij} - v_{ik}) \end{aligned}$$

where $F_{\varepsilon_{ik} - \varepsilon_{ij}}$ is the cdf of $\varepsilon_{ik} - \varepsilon_{ij}$.

- Different specifications of $F_{\varepsilon_{ik} - \varepsilon_{ij}}$ imply different models of multinomial response
- Suppose $j=1$ is the base category. Let $\bar{\varepsilon}_{ij} \equiv \varepsilon_{ij} - \varepsilon_{i1}$ and $\bar{v}_{ij} \equiv v_{ij} - v_{i1}$. Then

$$\begin{aligned} P[Y_i = j] &= P[\bar{\varepsilon}_{i2} \leq -\bar{v}_{i2}, \dots, \bar{\varepsilon}_{im} \leq -\bar{v}_{im}] \\ &= \int_{-\infty}^{-\bar{v}_{i2}} \dots \int_{-\infty}^{-\bar{v}_{im}} f_{\bar{\varepsilon}_{i2} \dots \bar{\varepsilon}_{im}}(t_2, \dots, t_m) dt_m \dots dt_2 \end{aligned}$$

which generally has no closed form solution.

In addition, expect that $\varepsilon_{i1}, \dots, \varepsilon_{im}$ to be correlated in general across alternatives, particularly amongst complements or substitutes.

On the other hand, easiest to start with the case where $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are iid. In particular, suppose $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are iid with a common type I extreme value distribution, i.e. continuous with common density

$$f_{\varepsilon_{ij}}(t) = e^{-t} e^{-e^{-t}}, \quad \text{for each } j$$

(In the case $m=2$, this implies a logit model) For general m , if $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are iid type I extreme value then

$$P[Y_i = j | V_{i1}, \dots, V_{im}] = \frac{v_{ij}}{\sum_{k=1}^m e^{v_{ik}}}$$

which is CL if $v_{ij} = \underline{X}_{ij}^T \underline{\beta}_0$

or MNL if $v_{ij} = \underline{X}_i^T \underline{\beta}_{0j}$

Suppose consumer i picks one of m alternatives. Let

$$U_{ij} = V_{ij} + \varepsilon_{ij}, \quad \text{for consumer } i, \quad \text{alternative } j \in \{1, \dots, m\}$$

(i.e. an ARUM) Let Y_i be the corresponding multinomial response variable. Then

$$\begin{aligned} P[i \text{ chooses } j] &= P[Y_i = j] \\ &= P[U_{ij} \geq U_{ik}, \quad \forall k \neq j] \\ &= P[V_{ij} + \varepsilon_{ij} \geq V_{ik} + \varepsilon_{ik}, \quad \forall k \neq j] \\ &= P[\varepsilon_{ik} - \varepsilon_{ij} \leq V_{ij} - V_{ik}, \quad \forall k \neq j] \end{aligned}$$

Let $F_{\varepsilon_{ik}-\varepsilon_{ij}}$ denote the cdf of $\varepsilon_{ik}-\varepsilon_{ij}$. Different specifications of $F_{\varepsilon_{ik}-\varepsilon_{ij}}$ imply different models of multinomial choice. Let $j=1$ be the base category. Let

$$\begin{aligned} \bar{\varepsilon}_{ij} &\equiv \varepsilon_{ij} - \varepsilon_{i1}, \quad j \in \{2, \dots, m\} \\ \bar{V}_{ij} &\equiv V_{ij} - V_{i1}, \quad j \in \{2, \dots, m\} \end{aligned}$$

Then

$$P[Y_i = 1] = P[\bar{\varepsilon}_{i2} \leq -\bar{V}_{i2}, \dots, \bar{\varepsilon}_{im} \leq -\bar{V}_{im}] = \int_{-\infty}^{-\bar{V}_{i2}} \dots \int_{-\infty}^{-\bar{V}_{im}} f_{\bar{\varepsilon}_{i2}, \dots, \bar{\varepsilon}_{im}}(t_2, \dots, t_m) dt_m \dots dt_2$$

Where $f_{\bar{\varepsilon}_{i2}, \dots, \bar{\varepsilon}_{im}}(\cdot)$ is the joint density of $\bar{\varepsilon}_{i2}, \dots, \bar{\varepsilon}_{im}$ (assuming absolute continuity of the corresponding joint distribution).

Note: $P[Y_i=j]$ has generally no closed form solution.

Easiest special case: Suppose $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are iid. In particular, suppose $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are iid type-I Extreme-value, with common density

$$f_{\varepsilon_{ij}}(t) = e^{-t}e^{-e^{-t}}, \quad \text{for each } j$$

(for $m=2$, this is a logit model) For general m , we get

$$\begin{aligned} P[Y_i = j | V_{i1}, \dots, V_{im}] &= \frac{V_{ij}}{\sum_{k=1}^m e^{V_{ik}}} \\ &\text{i. e. CL if } V_{ij} = \underline{X}_{ij}^T \underline{\beta}_0 \\ &\text{i. e. MNL if } V_{ij} = \underline{X}_i^T \underline{\beta}_{0j} \end{aligned}$$

Naturally, assuming independence of $\varepsilon_{i1}, \dots, \varepsilon_{im}$ is implausible if any two alternatives are similar/complementary/substitutes.

e.g. Suppose alternatives 2, 3 are similar to most consumers. Then if for a representative consumer i , ε_{i2} is large and negative, then U_{i2} will tend to be overpredicted. So U_{i3} will also tend to be overpredicted. i.e. ε_{i3} is also large and negative (this shows that $\varepsilon_{i2}, \varepsilon_{i3}$ are positively correlated).

Note: The red-bus blue-bus problem exists because in many cases it is unrealistic to assume independence of $\varepsilon_{i1}, \dots, \varepsilon_{im}$.

Next: It is therefore natural to assume that

$$\underline{\varepsilon}_i \equiv [\varepsilon_{i1} \quad \cdots \quad \varepsilon_{im}]^T \sim N_m(\underline{0}, \underline{\Sigma})$$

with $\underline{\Sigma}$ an arbitrary positive definite matrix.

This implies a multinomial probit (MNP) model (we will return to this shortly).

Welfare Analysis with an ARUM

- The idea is to place dollar values on the effect of changing one or more components of V_{ij}
- Suppose for consumer i the deterministic utility component is

$$V_{ij} = V(M_i - P_j, \underline{X}_{ij})$$

where M_i is income, P_j is the price of alternative j and \underline{X}_{ij} are other observables. Then

$$U_{ij} = V(M_i - P_j, \underline{X}_{ij}) + \varepsilon_{ij}$$

Suppose \underline{X}_{ij} changes from $\underline{X}_{ij}^{(1)}$ to $\underline{X}_{ij}^{(2)}$. The compensating variation (CV) is the resulting change in income required to U_{ij} at its initial level. So the highest utility possible with income M_i and characteristic $\underline{X}_{ij}^{(1)}$ is the same as that possible with income $M_i - CV_i$ and characteristics $\underline{X}_{ij}^{(2)}$, i.e.

$$\max_{j \in \{1, \dots, m\}} V(M_i - P_j, \underline{X}_{ij}^{(1)}) + \varepsilon_{ij} = \max_{j \in \{1, \dots, m\}} V(M_i - CV_i - P_j, \underline{X}_{ij}^{(2)}) + \varepsilon_{ij}$$

e.g. Suppose $m=2$ and

$$U_{ij} = M_i + X_{ij} + \varepsilon_{ij}, \quad j \in \{1, 2\}$$

Then if X_{ij} goes from $X_{ij}^{(1)}$ to $X_{ij}^{(2)}$, we have

- If alternative 1 is chosen before and after, i.e. if $U_{i1}^{(2)} \equiv M_i - CV_i + X_{i1}^{(2)} + \varepsilon_{i1} = M_i + X_{i1}^{(1)} + \varepsilon_{i1}$, then

$$CV_i = X_{i1}^{(2)} - X_{i1}^{(1)}$$

- Similarly, if alternative 2 is chosen before and after, then $CV_i = X_{i2}^{(2)} - X_{i1}^{(1)}$

$$CV_i = X_{i2}^{(2)} - X_{i1}^{(1)}$$

- If alternative 1 is chosen before and alternative 2 after, then

$$CV_i = X_{i2}^{(2)} - X_{i1}^{(1)} + \varepsilon_{i2} - \varepsilon_{i1}, \quad \text{i. e. } CV_i \text{ is random}$$

- If alternative 2 is chosen before and alternative 1 after, then

$$CV_i = X_{i1}^{(2)} - X_{i2}^{(1)} + \varepsilon_{i1} - \varepsilon_{i2}, \quad \text{i. e. } CV_i \text{ is also random}$$

Of interest is $E[CV_i]$, which involves integrating $\varepsilon_{i1}, \varepsilon_{i2}$. Now, go back to the problem of replacing the assumption that $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are independent. The easiest generalization involves nested logit. Start by assuming that $\varepsilon_i \equiv [\varepsilon_{i1} \ \dots \ \varepsilon_{im}]^T$ has a generalized extreme value (GEV) distribution (McFadden, 1978) with joint pdf

$$F_{\varepsilon}(\varepsilon_1, \dots, \varepsilon_m) = e^{-G(e^{-\varepsilon_1}, \dots, e^{-\varepsilon_m})}$$

where the function $G(t_1, \dots, t_m)$ is specified to make F_{ε} a proper cdf and the marginal distributions of $\varepsilon_1, \dots, \varepsilon_m$ to be well defined.

Note: Assuming an ARUM $U_{ij} = V_{ij} + \varepsilon_{ij}$. We have

$$P[Y_i = j | X_i] = e^{V_{ij}} \frac{G_j(e^{-V_{i1}}, \dots, e^{-V_{im}})}{G(e^{-V_{i1}}, \dots, e^{-V_{im}})}$$

where

$$G_j(t_1, \dots, t_m) = \frac{\partial G(t_1, \dots, t_m)}{\partial t_j}$$

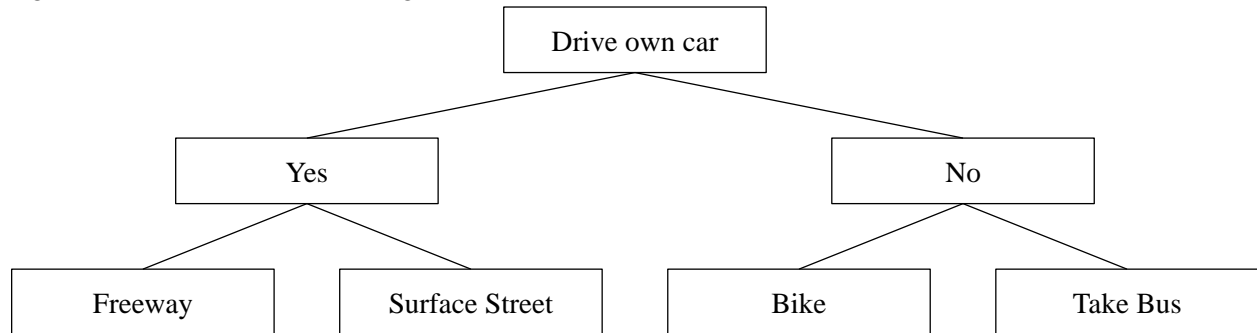
(Different specifications of $G(t_1, \dots, t_m)$ imply different models of multinomial response, e.g.

$$G(t_1, \dots, t_m) = \sum_{j=1}^m t_j$$

implies MNL)

The nested logit model involves a somewhat complicated specification of $G(\cdot)$. Nested logit assumes a multi-stage decision process.

e.g. Choice of mode for a morning commute



More generally, assume: The consumer has initially a mean of J choices (i.e. limbs). The j -th limb ($j=1, \dots, J$) has K_j branches labelled $(j1, \dots, jk, \dots, jK_j)$. The utility of the alternative attained by picking limb j and branch k (i.e. address (j,k)) is

$$U_{jk} = V_{jk} + \varepsilon_{jk}, \quad j \in \{1, \dots, J\}, \quad k \in \{1, \dots, K_j\}, \quad \text{where } \sum_{j=1}^J K_j = m$$

(this is a two-level decision tree; an arbitrary number of levels can be accommodated)

In a two-level model, we have

$$\begin{aligned} P[\text{individual } i \text{ chooses } (j, k)] &= P[\text{individual } i \text{ chooses limb } j] \\ &\quad \times P[\text{individual } i \text{ chooses branch } k | \text{individual } i \text{ chooses limb } j] \\ \Leftrightarrow P_{ijk} &= P_{ij} \times P_{ik|j} \end{aligned}$$

Nested logit involves specifying the random utility component $\underline{\varepsilon}$ in an ARUM as a generalized extreme value (GEV) random variable with cdf

$$F_{\underline{\varepsilon}}(\underline{\varepsilon}') = e^{-G\left(e^{-\varepsilon'_{i1}}, \dots, e^{-\varepsilon'_{iK_1}}, \dots, e^{-\varepsilon'_{j1}}, \dots, e^{-\varepsilon'_{jK_j}}\right)}$$

where

$$G\left(t_{11}, \dots, t_{1K_1}, \dots, t_{j1}, \dots, t_{jK_j}\right) = \sum_{j=1}^J \left(\sum_{k=1}^{K_j} Y_{jk}^{1/\rho_j} \right)^{\rho_j}, \quad \dots (*)$$

where the scale parameters are given by

$$\rho_j = \sqrt{1 - \text{Corr}(\varepsilon_{ijk}, \varepsilon_{ij\ell})}, \quad \text{for any } k, \ell \in \{1, \dots, K_j\}$$

Note: $\rho_j \equiv 1$ for each $j \in \{1, \dots, J\}$ iff the random utility components are independent (i.e. MNL)

Now: Let $Y_{ijk} \equiv 1 \{\text{individual } i \text{ chooses } (j, k)\}$. Then

$$P[Y_{ijk} = 1] = P[U_{ijk} \geq U_{i\ell m}, \quad \forall (\ell, m)]$$

Also: Assume that the deterministic utility component for alternative (j,k) is

$$V_{ijk} = \underline{Z}_{ij}^T \underline{\alpha}_0 + \underline{X}_{ijk}^T \underline{\beta}_{0j}, \quad j \in \{1, \dots, J\}, \quad k \in \{1, \dots, K_j\}$$

where \underline{Z}_{ij} varies over limbs but not branches. \underline{X}_{ijk} varies over limbs and branches. Then

$$P_{ijk} = P_{ij} \times P_{ik|j} = \frac{e^{\underline{Z}_{ij}^T \underline{\alpha}_0 + \rho_j I_{ij}}}{\sum_{m=1}^J e^{\underline{Z}_{im}^T \underline{\alpha}_0 + \rho_m I_{im}}} \times \frac{e^{\frac{1}{\rho_j} \underline{X}_{ijk}^T \underline{\beta}_j}}{\sum_{\ell=1}^{K_j} e^{\frac{1}{\rho_j} \underline{X}_{ij\ell}^T \underline{\beta}_j}}$$

where

$$I_{ij} = \log \left(\sum_{\ell=1}^{K_j} e^{\frac{1}{\rho_j} \underline{X}_{ij\ell}^T \underline{\beta}_j} \right)$$

is the inclusive value or log sum.

Warning: $\text{Corr}(\varepsilon_{ij1}, \varepsilon_{ij2}) = \text{Corr}(\varepsilon_{ij1}, \varepsilon_{ij3}) = \dots = \text{Corr}(\varepsilon_{ijK_j-1}, \varepsilon_{ijK_j}) = \rho_j \quad \forall j$ (branches share identical correlation)

Nested Logit

- Two level decision process
- J limbs, j-th limb has K_j branches

$$P[\text{individual } i \text{ chooses limb } j \text{ branch } k] = P[i \text{ chooses limb } j]P[i \text{ chooses branch } k | i \text{ chooses limb } j] \\ \Leftrightarrow P_{ijk} = P_{ij}P_{ik|j}$$

Assume that the deterministic utility component for alternative (j,k) is

$$V_{ijk} = \underline{Z}_{ij}^T \underline{\alpha}_0 + \underline{X}_{ijk}^T \underline{\beta}_{0j}, \quad j \in \{1, \dots, J\}, \quad k \in \{1, \dots, K_j\}$$

The nested logit model specifies

$$P_{ijk} = P_{ij}P_{ik|j} = \frac{e^{\underline{Z}_{ij}^T \underline{\alpha}_0 + \rho_{0j} I_{0ij}}}{\sum_{m=1}^J e^{\underline{Z}_{im}^T \underline{\alpha}_0 + \rho_{0m} I_{0im}}} \times \frac{e^{\frac{1}{\rho_{0j}} \underline{X}_{ijk}^T \underline{\beta}_{0j}}}{\sum_{\ell=1}^{K_j} e^{\frac{1}{\rho_{0j}} \underline{X}_{ij\ell}^T \underline{\beta}_{0j}}}$$

where

$$\rho_{0j} = \sqrt{1 - \frac{\text{Corr}[\varepsilon_{ijk}, \varepsilon_{ij\ell}]}{\text{assumed to be the same } \forall k \neq \ell}}, \quad \text{for any } k, \ell \in \{1, \dots, K_j\} \\ \text{and } I_{0ij} = \log \left(\sum_{\ell=1}^{K_j} e^{\frac{1}{\rho_{0j}} \underline{X}_{ij\ell}^T \underline{\beta}_{0j}} \right)$$

Let

$$\underline{\rho}_0 \equiv [\rho_{01} \quad \dots \quad \rho_{0J}]^T$$

How to estimate $\underline{\alpha}_0$, $\underline{\beta}_0$ and $\underline{\rho}_0$?

For individual I , we observe $m=K_1+\dots+K_J$ outcomes Y_{ijk} , where $Y_{ijk}=1 \{ \text{individual } i \text{ chooses } (j,k) \}$.

$$\text{Let } P_{ijk} \equiv P[Y_{ijk} = 1], \quad Y_{ij} \equiv \sum_{k=1}^{K_j} Y_{ijk}, \quad P_{ij} \equiv P[Y_{ij} = 1], \quad P_{ik|j} \equiv P[Y_{ijk} = 1 | Y_{ij} = 1]$$

Then $P_{ijk}=P_{ik|j}P_{ij}$ and the density of

$$\underline{Y}_i \equiv [Y_{i11} \quad \dots \quad Y_{i1K_2} \quad \dots \quad Y_{iJ1} \quad \dots \quad Y_{iJK_J}]$$

is given by

$$f_Y(\underline{y}_i; \underline{\alpha}_0, \underline{\beta}_0, \underline{\rho}_0) = \prod_{j=1}^J \prod_{k=1}^{K_j} (P_{ik|j}P_{ij})^{y_{ijk}} = \prod_{j=2}^J p_{ij}^{y_{ij}} \prod_{k=1}^{K_j} P_{ik|j}^{y_{ijk}}$$

Therefore the joint density of the sample $\underline{Y}_1, \dots, \underline{Y}_n$ is

$$\prod_{i=1}^n f_Y(\underline{y}_i; \underline{\alpha}_0, \underline{\beta}_0, \underline{\rho}_0)$$

and the full information maximum-likelihood (FIML) estimator maximizes

$$\frac{1}{n} \log \hat{L}_n(\underline{\alpha}, \underline{\beta}, \underline{\rho}) \equiv \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^J Y_{ij} \log P_{ij} + \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^J \sum_{k=1}^{K_j} Y_{ijk} \log P_{iklj}, \quad \dots (*)$$

with respect to $\underline{\alpha}, \underline{\beta}, \underline{\rho}$.

Important: In general, $\frac{1}{n} \log \hat{L}_n(\underline{\alpha}, \underline{\beta}, \underline{\rho})$ is not globally concave \rightarrow May be useful to use limited information maximum-likelihood (LIML) to generate starting values.

LIML in this context first estimates the parameters in the second term of the RHS of (*), which is the log-likelihood of a CL model with parameter $\frac{1}{\rho_j} \underline{\beta}_j$ for each limb. LIML then estimates the parameters in the first term of the RHS of (*), which is the log-likelihood of another CL model with added regressors \hat{I}_{ij} which embed the estimates from the first-stage estimation of the second term of (*).

Some discussion regarding the nested logit model

- (Obvious) Not all choice problems have a natural nesting structure. In practice, may need to treat the nesting structure as another parameter to estimate, perhaps via AIC/BIC
- Can occasionally get ML estimates of the scale parameters ρ_j that lie outside the range $[0,1]$, which implies a violation of the ARUM specification
- Often see decision trees with more than two levels

e.g. Consumer behavior in a car dealership

1. Buy a car (Y/N)
2. If Y, buy a new car (Y/N)
3. If Y, which class of new car to buy (e.g. sedan, minivan, SUV, etc.)
4. Given the choice at level 3, foreign/domestic
5. Given the choice at level 4, which model to buy

Next: Consider random parameters logit (RPL)

- Another generalization of CL or MNL in which the random utility components can exhibit correlation across alternatives
- For individual i , suppose the utility of alternative j is

$$U_{ij} = \underline{X}_{ij}^T \underline{\beta}_i + \varepsilon_{ij}, \quad j \in \{1, \dots, m\}$$

where $\varepsilon_{i1}, \dots, \varepsilon_{im}$ are iid type-I extreme value (as in CL), but where $\underline{\beta}_i \sim N(\underline{\beta}_0, \underline{\Sigma}_{\beta_0})$, i.e. random coefficients. Assume that $\underline{\beta}_i$ is independent of ε_{ij} for all j .

Therefore

$$U_{ij} = \underline{X}_{ij}^T \underline{\beta}_0 + \underline{X}_{ij}^T \underline{u}_i + \varepsilon_{ij}, \quad \dots (\dagger)$$

where $\underline{u}_i \sim N(\underline{0}, \underline{\Sigma}_{\underline{\beta}_0})$

Let

$$\zeta_{ij} \equiv \underline{X}_{ij}^T \underline{u}_i + \varepsilon_{ij}, \quad \text{then } \text{Cov}[\zeta_{ij}, \zeta_{ik}] = \underline{X}_{ij}^T \underline{\Sigma}_{\underline{\beta}_0} \underline{X}_{ik} \neq 0, \quad \text{for } j \neq k \text{ in general}$$

In most RPL applications, $\underline{\Sigma}_{\underline{\beta}_0}$ is diagonal with some diagonal elements equal to zero, i.e. $\underline{\beta}_i$ is singular Normal (often see (\dagger) aggregated over I to model market demand)

How to estimate RPL models?

- Cannot apply CL to estimate $\underline{\beta}_0$ in (\dagger) ; inconsistent
- However, can integrate out $\underline{\beta}_i$ to get

$$P_{ij} \equiv P[Y_i = j] = \int \frac{e^{\underline{X}_{ij}^T \underline{\beta}_i}}{\sum_{k=1}^m e^{\underline{X}_{ik}^T \underline{\beta}_i}} \varphi(\underline{\beta}_i; \underline{\beta}_0, \underline{\Sigma}_{\underline{\beta}_0}) d\underline{\beta}_i$$

a multivariate integral, where $\varphi(\underline{\beta}_i; \underline{\beta}_0, \underline{\Sigma}_{\underline{\beta}_0})$ is the $N(\underline{\beta}_0, \underline{\Sigma}_{\underline{\beta}_0})$ -density.

The MLE of $\underline{\beta}_0, \underline{\Sigma}_{\underline{\beta}_0}$ maximizes

$$\frac{1}{n} \log \hat{L}_n(\cdot) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m Y_{ij} \log P_{ij}, \quad \text{wrt } \underline{\beta}, \underline{\Sigma}_{\underline{\beta}}$$

Major problem: P_{ij} has no closed form, so P_{ij} needs to be approximated, usually by simulation.

A popular method of approximating P_{ij} : Take the average of S evaluations of the integrand at each S random draws of $\underline{\beta}_i$ from a $N(\underline{\beta}_0, \underline{\Sigma}_{\underline{\beta}_0})$ -distribution ($\underline{\beta}_0, \underline{\Sigma}_{\underline{\beta}_0}$ are unknown, so the simulation is part of an iterative procedure that produces $\underline{\beta}_0^{(r)}$ and $\underline{\Sigma}_{\underline{\beta}_0}^{(r)}$ at the iteration)

The maximum simulated likelihood (MSL) estimator maximizes

$$\frac{1}{n} \log \hat{L}_n(\underline{\beta}, \underline{\Sigma}_{\underline{\beta}}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m Y_{ij} \log \left(\frac{1}{S} \sum_{s=1}^S \frac{e^{\underline{X}_{ij}^T \underline{\beta}_i^{(s)}}}{\sum_{k=1}^m e^{\underline{X}_{ik}^T \underline{\beta}_i^{(s)}}} \right)$$

where $\underline{\beta}_i^{(s)}$ ($s=1, \dots, S$) are iid $N(\underline{\beta}_0^{(r)}, \underline{\Sigma}_{\underline{\beta}_0}^{(r)})$.

Consistency requires $n \rightarrow \infty, S \rightarrow \infty$ with $\frac{\sqrt{n}}{S} \rightarrow 0$: Gouriéroux and Monfort (1991, *Annales d'Économie et de Statistique*).

Christian Gouriéroux and Alain Monfort (1991), Simulation based inference in models with heterogeneity, *Annales d'Économie et de Statistique* (20/21), 69–107

Extension: Can impose random parameters on nested logit models.

Consider a generalized random utility model, where individual i and alternative j

$$U_{ij} = V_{ij}(\underline{X}_i, \underline{\beta}_0) + \varepsilon_{ij}$$

where \underline{X}_i is a vector of observed regressors, $\underline{\beta}_0$ is an unknown vector of constants and $\{\varepsilon_{ij}\}$ are independent over i but (in general) correlated over j .

Assume that the distribution of ε_{ij} is such that a closed form for the choice probabilities exists, with

$$P_{ij} \equiv P[\text{individual } i \text{ chooses alternative } j] = G_j(\underline{V}_i(\underline{X}_i, \underline{\beta}_0), \underline{\theta}_\varepsilon)$$

where $\underline{V}_i(\underline{X}_i, \underline{\beta}_0) = [V_{i1}(\underline{X}_i, \underline{\beta}_0) \quad \cdots \quad V_{im}(\underline{X}_i, \underline{\beta}_0)]^T$ and where $\underline{\theta}_\varepsilon$ is a vector of unknown parameters of the distribution of

$$\underline{\varepsilon}_i \equiv [\varepsilon_{i1} \quad \cdots \quad \varepsilon_{im}]^T$$

(e.g. $\underline{\varepsilon}_i$ could have a GEV distribution)

Can introduce additional randomness to this model by specifying

$$V_{ij} = V_{ij}(\underline{X}_i, \underline{\xi}_i, \underline{\beta}_0)$$

where $\underline{\xi}_i$ is a random vector.

In this case,

$$P_{ij} = \int G_j(\underline{V}_i(\underline{X}_i, \underline{\xi}_i, \underline{\beta}_0), \underline{\theta}_\varepsilon) f(\underline{\xi}_i; \underline{\theta}_\xi) d\underline{\xi}_i$$

where $f(\underline{\xi}; \underline{\theta}_\xi)$ is the density of $\underline{\xi}$.

Can motivate the introduction of additional randomness in this case by assuming that individuals belong to one of C latent classes; suppose that $\underline{\beta}$, $\underline{\theta}_\varepsilon$ vary by class. Then

$$P_{ij} = \sum_{c=1}^C G_j(\underline{V}_i(\underline{X}_i, \underline{\beta}^c), \underline{\theta}_\varepsilon^c) \pi_c$$

where $\pi_c \in [0,1]$, and $\sum_c \pi_c = 1$: Walker and Ben-Akiva (2002, *Mathematical Social Sciences*)

Joan Walker and Moshe Ben-Akiva (2002), Generalized random utility model, *Mathematical Social Sciences* 43 (3), 303–343

Generalized random utility models

Suppose that individuals belong to one of C latent classes. If the utility parameters vary by class

$$P_{ij} = \sum_{c=1}^C \int G_j \left(V_i(\underline{X}_i, \underline{\xi}_i, \underline{\beta}^c), \underline{\theta}_{\xi}^c \right) f(\underline{\xi}_i; \underline{\theta}_{\xi}) d\underline{\xi}_i \pi_c$$

where $\pi_c \in [0,1]$ with $\sum_c \pi_c = 1$ (usually assume $c \in \{2,3\}$).

The corresponding MSL estimator maximizes

$$\frac{1}{n} \log \hat{L}_n(\underline{\beta}, \underline{\Sigma}_{\underline{\beta}}) \equiv \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m Y_{ij} \log \left(\frac{1}{S} \sum_{s=1}^S \sum_{c=1}^C G_j \left(V_i(\underline{X}_i, \underline{\xi}_i^s, \underline{\beta}^c), \underline{\theta}_{\xi}^c \right) \pi_c \right)$$

where $\underline{\xi}_i^s$ ($s \in \{1, \dots, S\}$) are iid with density $f(\cdot | \underline{\theta}_{\xi})$ (Walker and Ben-Akiva, 2002, *Mathematical Social Sciences*)

Lastly: Consider an obvious way of modelling correlation amongst the elements of

$$\underline{\varepsilon}_i = [\varepsilon_{i1} \quad \dots \quad \varepsilon_{im}]^T$$

The multinomial probit (MNP) model assumes for individual i and alternative j that

$$U_{ij} = V_{ij} + \varepsilon_{ij}, \quad \text{where } \underline{\varepsilon}_i = [\varepsilon_{i1} \quad \dots \quad \varepsilon_{im}]^T \sim N_m(\underline{0}, \underline{\Sigma})$$

where usually $V_{ij} = \underline{X}_{ij}^T \underline{\beta}$ or $V_{ij} = \underline{X}_i^T \underline{\beta}_j$. Need to impose restrictions on $\underline{\Sigma}$ to guarantee identification.

Recall

$$P_{ij} \equiv P[Y_i = j] = P[\varepsilon_{ik} - \varepsilon_{ij} \leq V_{ij} - V_{ik}, \quad \forall k \neq j]$$

Normalize the utility differences w.r.t. the utility of alternative 1; we can then impose $\varepsilon_{i1} \equiv 0 \forall i$ and restrict the (1,1)-element of $\underline{\Sigma}$ to zero as well as several other elements of $\underline{\Sigma}$.

e.g. If $m=2$, and we impose $\varepsilon_{i1} \equiv 0$, then $\sigma_{11} = \sigma_{12} = \sigma_{21} = 0$. If also impose $\sigma_{22} = 1$, then

$$\varepsilon_{i2} - \varepsilon_{i1} = \varepsilon_{i2} \sim N(0,1)$$

which implies a probit model of binary response. One often in the context of MNP models sees factor models of the form

$$\varepsilon_{ij} = \zeta_{ij} + \sum_{\ell=1}^L c_{ij\ell} \xi_{i\ell}, \quad j \in \{1, \dots, m\}$$

where ζ_{ij} , ξ_{i1} , \dots , ξ_{iL} are iid $N(0,1)$ and $c_{ij\ell}$ are factor loadings that can be estimated.

In this case the number of distinct parameters in $\underline{\Sigma}$ to estimate goes from $\frac{m(m+1)}{2}$ to L , which is invariably smaller.

How to estimate MNP models?

No closed form for the choice probabilities.

e.g. If $m=3$ with $\varepsilon_{i1} \equiv 0$, then

$$P_{ij} = P[Y_i = j] \int_{-\infty}^{-\bar{V}_{i31}} \int_{-\infty}^{-\bar{V}_{i21}} f(\bar{\varepsilon}_{i21}, \bar{\varepsilon}_{i31}) d\bar{\varepsilon}_{i21} d\bar{\varepsilon}_{i31}$$

where

$$\bar{\varepsilon}_{i21} \equiv \varepsilon_{i2} - \varepsilon_{i1} = \varepsilon_{i2}, \quad \bar{\varepsilon}_{i31} \equiv \varepsilon_{i3} - \varepsilon_{i1} = \varepsilon_{i3}$$

and where $f(\bar{\varepsilon}_{i21}, \bar{\varepsilon}_{i31})$ is a (singular) bivariate Normal density with up to two free parameters in the covariance matrix.

In general, a scenario with m choices requires the evaluation of an $(m-1)$ -fold integral, which is usually impossible when $m > 4$.

When m is large, usually use simulation to approximate P_{ij} and estimate $\underline{\beta}_0, \underline{\Sigma}_0$ by MSL, i.e. by maximizing

$$\frac{1}{n} \log \hat{L}_n(\underline{\beta}, \underline{\Sigma}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m Y_{ij} \log \hat{P}_{ij}$$

where \hat{P}_{ij} is obtained by simulation.

A logical alternative to MSL (especially in the context of MNP) is to use the method of simulated moments (MSM), i.e. to solve

$$\sum_{i=1}^n \sum_{j=1}^m (Y_{ij} - \hat{P}_{ij}) Z_i = 0$$

where \hat{P}_{ij} is obtained by simulation such that

$$(Y_{ij} - \hat{P}_{ij}) Z_i \text{ is unbiased for } (Y_{ij} - P_{ij}) Z_i$$

Some discussion

- MNP generally involves a more challenging approximation to P_{ij} (via simulation) than is the case with MSL applied to RPL
- This likely explains why nested logit and RPL tend to be more popular than MNP

Example: Choice of fishing mode (based on a sample of $n=1,182$ recreational fishers; Thompson and Crooke (1991))

For individual $i \in \{1, \dots, 1,182\}$, observe $Y_i \in \{1, 2, 3, 4\}$ according to whether i chooses to fish at a beach ($j=1$), off a pier ($j=2$), off a private boat ($j=3$) or off a charter boat ($j=4$). Expect that choice of mode depends on relative prices of each mode and the catch rate for each mode.

Cynthia J. Thompson and Stephen J. Crooke (1991), Results of the southern California sportfish economic survey, Southwest Fisheries Center

Variable	In-sample averages			
	Y=1 (beach)	Y=2 (pier)	Y=3 (private)	Y=4 (charter)
P _{beach} (\$)	36	31	138	121
P _{pier} (\$)	36	31	138	121
P _{private} (\$)	98	82	42	45
P _{charter} (\$)	125	110	71	75
CR _{beach} (%)	28	26	21	25
CR _{pier} (%)	22	20	13	16
CR _{private} (%)	16	15	18	18
CR _{charter} (%)	52	50	65	69
Sample proportions (%)	11.3	15.1	35.4	38.2
Observations	134	178	418	452

So

- Fishers tend to pick the mode where it is cheapest for them to fish
- Beach and pier fishers find it much cheaper to fish from shore than from a boat
- Relationship between catch rates and mode choice is not obvious; charter boats consistently have the highest catch rate

Clearly have alternative-varying regressors, so a conditional logit (CL) is natural

$$P_{ij} = P[Y_i = j] = \frac{e^{\beta_p P_{ij} + \beta_c CR_{ij}}}{\sum_{k=1}^4 e^{\beta_p P_{ik} + \beta_c CR_{ik}}}$$

We find $\hat{\beta}_p = -.021$ and $\hat{\beta}_c = .953$. So decrease in the price of one alternative increases the probability of choosing that alternative and decreases the probability of every other alternative.

Similar story if catch rates change. Customary to report estimated average marginal effects, i.e. estimates of

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial P_{ij}}{\partial X_{ikr}}$$

i.e. the average marginal response of the probability of choosing j when the r-th regressor for alternative k increases by a tiny amount, all else equal.

Given the CL specification, the average response is estimated as

$$\frac{1}{n} \sum_{i=1}^n \hat{P}_{ij} (\delta_{ijk} - \hat{P}_{ik}) \hat{\beta}_r = \bar{P}_j (\delta_{jk} - \bar{P}_k) \hat{\beta}_r$$

(where \hat{P}_{ij} are predicted choice probabilities, \bar{P}_j are in-sample proportions of each mode) Therefore: Can report the average effect of a \$100 increase in the price of each mode, all else constant.

Average marginal response				
	Y=1	Y=2	Y=3	Y=4
Variable	(beach)	(pier)	(private)	(charter)
$\Delta P[\text{beach}]$ (%)	-27.2	11.9	8.5	6.8
$\Delta P[\text{pier}]$ (%)	11.9	-26.3	8.0	6.4
$\Delta P[\text{private}]$ (%)	8.0	8.0	-39.1	22.5
$\Delta P[\text{charter}]$ (%)	6.8	6.4	22.5	-35.7

We see that shore-based and boat-based modes have a high degree of substitutability. Now suppose that individual choice of fishing mode also depends on monthly income (in thousands of dollars).

Average incomes across the four modes					
	Y=1	Y=2	Y=3	Y=4	
Variable	(beach)	(pier)	(private)	(charter)	(overall)
Income (\$1K/month)	4.052	3.387	4.654	3.881	4.099

Income is obviously invariant across alternatives, so a multinomial logit (MNL) model can be specified

$$P_{ij} = P[Y_i = j] = \frac{e^{\alpha_j + \beta_{ij}I_i}}{\sum_{k=1}^4 e^{\alpha_k + \beta_{ik}I_i}}, \quad j \in \{1, \dots, 4\}$$

Can normalize $\alpha_1 = \beta_{11} = 0$ to ensure that probabilities sum to one (this imposes beach fishing as the baseline category).

Example: Choice of fishing mode. Sample of $n=1,182$ recreational fishers (Thompson and Crooke, 1991).

For fisher i we have $Y_i \in \{1, \dots, 4\}$ according to which of four modes is chosen. Last week we looked at the effect of relative prices on choice of mode.

Now: Suppose that mode choice is also depends on monthly income (in thousands of dollars).

Subsample averages:					
	Y=1	Y=2	Y=3	Y=4	
Monthly income	(beach)	(pier)	(private)	(charter)	(overall)
(thousands of dollars)	4.052	3.387	4.654	3.881	4.099

Income doesn't vary across alternatives, so consider an MNL specification.

$$P_{ij} = \frac{e^{\alpha_j + \beta_{ij} I_i}}{\sum_{k=1}^4 e^{\alpha_k + \beta_{ik} I_i}}, \quad j \in \{1, \dots, 4\}$$

Normalize $\alpha_1=0$, $\beta_{11}=0$. So:

$\hat{\alpha}_1(\text{beach}) = 0.000$	$\hat{\beta}_{11}(\text{beach}) = 0.000$
$\hat{\alpha}_2(\text{pier}) = 0.814$	$\hat{\beta}_{12}(\text{pier}) = -0.143$
$\hat{\alpha}_3(\text{private}) = 0.739$	$\hat{\beta}_{13}(\text{private}) = 0.092$
$\hat{\alpha}_4(\text{charter}) = 1.341$	$\hat{\beta}_{14}(\text{charter}) = -0.032$

The average marginal effect of an income change is

$$\frac{1}{n} \sum_{i=1}^n \frac{\partial P_{ij}}{\partial I_i}, \quad \text{which is estimated as } \frac{1}{n} \sum_{i=1}^n \hat{P}_{ij} (\hat{\beta}_{ij} - \hat{\beta}_i), \quad \text{where } \hat{\beta}_i = \sum_{j=1}^4 \hat{P}_{ij} \hat{\beta}_{ij}$$

(Exercise: show this) Therefore the average effect of a \$1,000 increase in monthly income is

$$\Delta P[\text{beach}] = 0.000 \quad \Delta P[\text{pier}] = -0.021 \quad \Delta P[\text{private}] = 0.033 \quad \Delta P[\text{charter}] = -0.012$$

So: A \$1,000 increase in monthly income tends on average to cause fishers to move out of pier and charter boat fishing into private boat fishing.

Goodness of fit:

The CL model has a pseudo- R^2 of 0.162
The MNL model has a pseudo- R^2 of 0.099

Can do better than CL or MNL by adopting mixed logit (i.e. use prices and monthly income as regressors)

$$P_{ij} = \frac{e^{\beta_P P_{ij} + \beta_C C_{ij} + \alpha_j + \beta_{Ij} I_i}}{\sum_{k=1}^4 e^{\beta_P P_{ik} + \beta_C C_{ik} + \alpha_k + \beta_{Ik} I_i}}$$

which is equivalent to the CL model

$$P_{ij} = \frac{e^{\beta_P P_{ij} + \beta_C C_{ij} + \sum_{\ell=1}^4 \alpha_{\ell} d_{ij\ell} + \beta_{I\ell} dI_{ij\ell}}}{\sum_{k=1}^4 e^{\beta_P P_{ik} + \beta_C C_{ik} + \sum_{\ell=1}^4 \alpha_{\ell} d_{ik\ell} + \beta_{I\ell} dI_{ik\ell}}}$$

where $d_{ij\ell} = 1\{j = \ell\}$ and $dI_{ij\ell} = d_{ij\ell} I_i = I_i 1\{j = \ell\}$ (exercise: verify this).

We find that $\hat{\beta}_P = -0.025$, $\hat{\beta}_C = 0.358$.

$\hat{\alpha}_1(\text{beach}) = 0.000$	$\hat{\beta}_{I1}(\text{beach}) = 0.000$
$\hat{\alpha}_2(\text{pier}) = 0.778$	$\hat{\beta}_{I2}(\text{pier}) = -0.128$
$\hat{\alpha}_3(\text{private}) = 0.527$	$\hat{\beta}_{I3}(\text{private}) = 0.089$
$\hat{\alpha}_4(\text{charter}) = 1.694$	$\hat{\beta}_{I4}(\text{charter}) = -0.033$

- Compared to the CL fit, $\hat{\beta}_C$ in the mixed logit is much smaller
- While the others are not much different from their counterparts in the CL or MNL models
- However: The pseudo- R^2 is 0.258, so much higher than either the CL pseudo- R^2 or the MNL pseudo- R^2

Models with censored or truncated dependent variables

Let Y^* denote a latent dependent variable, i.e. Y^* is not completely observed; only observe a function Y of Y^* . For example,

Censoring If censoring is from below we observe

$$Y = \max\{Y^*, L\} = \begin{cases} Y^*, & Y^* > L \\ U, & Y^* \leq L \end{cases}$$

If censoring is from above we observe

$$Y = \min\{Y^*, U\} = \begin{cases} Y^*, & Y^* < U \\ U, & Y^* \geq U \end{cases}$$

e.g. Only observed offered wages Y^* if Y^* exceeds a reservation wage L (censoring from below)

e.g. Incomes above U dollars per annum are reported as U , i.e. incomes above U are top coded (censoring from above)

Truncation If truncation is from below we observe $Y=Y^*$ iff $Y^* > L$; nothing is reported if $Y^* \leq L$; i.e. there is information loss. Similarly, if truncation is from above we observe $Y=Y^*$ if $Y^* < U$; nothing is reported when $Y^* \geq U$.

MLE: Censored MLE: Suppose censoring is from below: If $Y > L$, then the conditional density of Y given \underline{X} is the same as that of Y^* , i.e.

$$f_Y(y|\underline{X}) = f_{Y^*}(y|\underline{X})$$

If $Y=L$ then

$$P[Y = L|\underline{X}] = P[Y^* \leq L|\underline{X}] = F_Y(L|\underline{X})$$

(Therefore the conditional distribution of Y given \underline{X} is neither absolutely continuous nor discrete)

We have

$$f_Y(y|\underline{X}) = f_{Y^*}(y|\underline{X})^d F_{Y^*}(L|\underline{X})^{1-d}$$

where $d=1\{y>L\}$.

Now: Suppose Y^* has a conditional distribution given \underline{X} parameterized with $\underline{\theta} \in \Theta$. Then the log-likelihood for a sample $[\underline{X}_i^T \ Y_i \ D_i]$ ($i=1, \dots, n$) is

$$\hat{L}_n(\underline{\theta}) \equiv \sum_{i=1}^n [D_i \log f_{Y^*}(y_i|\underline{X}_i; \underline{\theta}) + (1 - D_i) \log F_{Y^*}(L|\underline{X}_i; \underline{\theta})], \quad \dots (*)$$

So: We see that ignoring censoring (i.e. using only observations for which $Y_i > L$) is equivalent to ignoring the second term of (*). This leads to an inconsistent estimate of $\underline{\theta}$.

Truncated MLE: Suppose truncation is from below. Then the conditional density of the observed Y given \underline{X} is

$$\begin{aligned} f_Y(y|\underline{X}) &= f_{Y^*}(y|Y^* > L, \underline{X}) \\ &= \frac{f_{Y^*}(y|\underline{X})}{P[Y^* > L|\underline{X}]} \\ &= \frac{f_{Y^*}(y|\underline{X})}{1 - F_{Y^*}(L|\underline{X})} \end{aligned}$$

Suppose Y^* has a conditional distribution given \underline{X} parameterized by $\underline{\theta} \in \Theta$. Then the log-likelihood for a sample $[\underline{X}_i^T \ Y_i]$ ($i=1, \dots, n$) is

$$\hat{L}_n(\underline{\theta}) = \sum_{i=1}^n (\log f_{Y^*}(Y_i|\underline{X}_i; \underline{\theta}) - \log(1 - F_{Y^*}(L|\underline{X}_i; \underline{\theta}))), \quad \dots (\dagger)$$

So ignoring truncation is equivalent to ignoring the second term of (\dagger), which leads to an inconsistent estimate of $\underline{\theta}$.

The most famous example of a regression with a censored dependent variable is the Tobit model (Tobin, 1958).

James Tobin (1958), Estimation of relationships for limited dependent variables, *Econometrica* 26 (1), 24–36

Assume that

$$Y^* = \underline{X}^\top \underline{\beta}_0 + \varepsilon, \quad \text{where } \varepsilon|\underline{X} \sim N(0, \sigma_0^2)$$

Therefore

$$Y^*|\underline{X} \sim N(\underline{X}^\top \underline{\beta}_0, \sigma_0^2)$$

Suppose we observe

$$Y = \begin{cases} Y^*, & Y^* > 0 \\ 0, & Y^* \leq 0 \end{cases}$$

We have

$$\begin{aligned} F_{Y^*}(0|\underline{X}) &= P[Y^* \leq 0|\underline{X}] = P[\underline{X}^\top \underline{\beta}_0 + \varepsilon \leq 0|\underline{X}] \\ &= \Phi\left(-\frac{\underline{X}^\top \underline{\beta}_0}{\sigma_0}\right) \\ &= 1 - \Phi\left(\frac{\underline{X}^\top \underline{\beta}_0}{\sigma_0}\right) \end{aligned}$$

So the conditional density of Y given \underline{X} is

$$f_Y(y|\underline{X}) = \left[\frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2\sigma_0^2}(y - \underline{X}^\top \underline{\beta}_0)^2} \right]^d \left[1 - \Phi\left(\frac{\underline{X}^\top \underline{\beta}_0}{\sigma_0}\right) \right]^{1-d}, \quad \text{where } d = 1\{y > 0\}$$

The Tobit MLE $\hat{\theta}_n \equiv [\hat{\beta}_n^\top \quad \hat{\sigma}_n^2]^\top$ maximizes

$$\hat{L}_n(\theta) = \sum_{i=1}^n \left\{ D_i \left[-\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (Y_i - \underline{X}_i^\top \underline{\beta})^2 \right] + (1 - D_i) \log \left(1 - \Phi\left(\frac{\underline{X}_i^\top \underline{\beta}}{\sigma}\right) \right) \right\}$$

which has partial derivatives (verify)

$$\begin{aligned} \frac{\partial \hat{L}_n(\theta)}{\partial \underline{\beta}^\top} &= \sum_{i=1}^n \left\{ \frac{1}{\sigma^2} \left[D_i \frac{Y_i - \underline{X}_i^\top \underline{\beta}}{\sigma^2} - (1 - D_i) \frac{\varphi\left(\frac{\underline{X}_i^\top \underline{\beta}}{\sigma}\right)}{1 - \Phi\left(\frac{\underline{X}_i^\top \underline{\beta}}{\sigma}\right)} \right] \underline{X}_i \right\} \\ \frac{\partial \hat{L}_n(\theta)}{\partial \sigma^2} &= \sum_{i=1}^n \left\{ D_i \left[-\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} (Y_i - \underline{X}_i^\top \underline{\beta})^2 \right] + (1 - D_i) \frac{\varphi\left(\frac{\underline{X}_i^\top \underline{\beta}}{\sigma}\right)}{1 - \Phi\left(\frac{\underline{X}_i^\top \underline{\beta}}{\sigma}\right)} \frac{\underline{X}_i^\top \underline{\beta}}{2\sigma^3} \right\} \end{aligned}$$

Exercise: Derive the limiting distribution as $n \rightarrow \infty$ of the Tobit MLE.

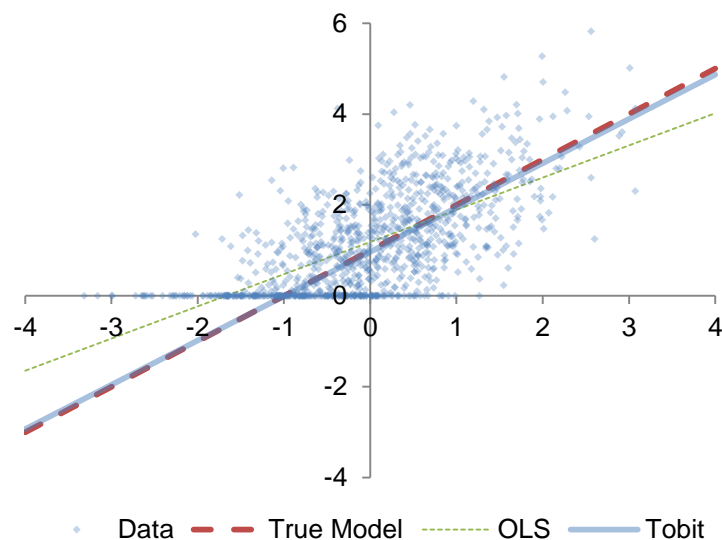
Fragility of the Tobit MLE: If the latent error ε is heteroskedastic or non-Normal then the Tobit MLE is inconsistent. To see this, observe that if the score identity holds, i.e. if

$$E \left[\frac{\partial \hat{L}_n(\theta)}{\partial \beta^\top} \bigg|_{\underline{\beta} = \underline{\beta}_0} \right] = \underline{0}$$

then $E[D_i] = \Phi \left(\frac{\underline{X}_i^\top \underline{\beta}_0}{\sigma_0} \right)$ and $E[D_i Y_i] = \Phi \left(\frac{\underline{X}_i^\top \underline{\beta}_0}{\sigma_0} \right) \underline{X}_i^\top \underline{\beta}_0 + \sigma_0 \varphi \left(\frac{\underline{X}_i^\top \underline{\beta}_0}{\sigma_0} \right)$. Under heteroskedasticity

$$E[D_i] = \Phi \left(\frac{\underline{X}_i^\top \underline{\beta}_0}{\sigma_i} \right) \neq \Phi \left(\frac{\underline{X}_i^\top \underline{\beta}_0}{\sigma_0} \right), \quad \text{unless } \sigma_i \equiv \sigma_0 \text{ for all } i$$

i.e. heteroskedasticity implies failure of the score identity.



True DGP ($n=1000, \alpha=\beta=\sigma^2=1$)
 OLS ($\hat{\alpha} = 1.19, \hat{\beta} = 0.71, \hat{\sigma}^2 = 0.73$)
 Tobit ($\hat{\alpha} = 0.97, \hat{\beta} = 0.98, \hat{\sigma}^2 = 1.07$)

Final Exam

- Cumulative
- Wednesday, December 21, 2016
- 10–12 AM
- In Bolton B80