

Maximum Likelihood Methods

Jesper Riis–Vestergaard Sørensen

University of Copenhagen, Department of Economics

Maximum Likelihood Estimation: Aim

- ▶ Previously: Modelled [feature(s)] of distribution of $y|\mathbf{x}$.
 - ▶ E.g. $E(y|\mathbf{x})$ and $\text{var}(y|\mathbf{x})$.
- ▶ Maximum likelihood estimation (MLE) more ambitious.
- ▶ Model for *entire distribution* $D(y|\mathbf{x})$

Why MLE? Advantages

Efficiency

- ▶ MLE uses entire $D(y|\mathbf{x})$.
- ▶ Structure \Rightarrow Information.

May estimate any feature

- ▶ Conditional moments: $E(y|\mathbf{x})$ and $\text{var}(y|\mathbf{x})$.
- ▶ Conditional prob's: $P(y = 1|\mathbf{x})$ and $P(y \in [a, b]|\mathbf{x})$.
- ▶ Conditional density.
- ▶ Derivatives (wrt. \mathbf{x}) thereof...

Why Not MLE? Drawbacks

Nonrobustness

- ▶ MLE uses entire $D(y|\mathbf{x}) \dots$
- ▶ Inconsistent (in general) if misspecified.
- ▶ (Exceptions exist.)

Outline

Framework

Example: Probit

Identification and Solution Uniqueness

Asymptotic Properties

Consistency

Asymptotic Normality

Asymptotic Variance Estimation

Example: Probit Avar Estimation

Framework

Truth vs. Model

Object of interest: “True” density $p_o(\mathbf{y}|\mathbf{x})$ of $\mathbf{y}_i|\mathbf{x}_i$.

- ▶ Possible values $(\mathbf{y}, \mathbf{x}) \in \mathcal{Y} \times \mathcal{X}$ for $(\mathbf{y}_i, \mathbf{x}_i)$.
- ▶ Discrete and/or continuous elements allowed.
- ▶ Only discrete: Integrals \rightarrow Sums.

Parametric model: $f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^P$.

Model Assumptions

Parametric model: $f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^P$.

Assume

1. Legitimate densities:

$$\begin{aligned} f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) &\geq 0, \text{ all } (\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}), \\ \int_{\mathcal{Y}} f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}) \nu(\mathrm{d}y) &= 1, \text{ all } (\mathbf{x}, \boldsymbol{\theta}). \end{aligned}$$

2. Correct specification: For *some* $\boldsymbol{\theta}_o \in \Theta$,

$$p_o(\mathbf{y}|\mathbf{x}) = f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}_o), \text{ all } (\mathbf{y}, \mathbf{x}).$$

Identification in Maximum Likelihood Context

Definition: θ_o **identified** if and only if for all $\theta \in \Theta \setminus \{\theta_o\}$ s.t.

$$f(\mathbf{y}|\mathbf{x}; \theta) \neq f(\mathbf{y}|\mathbf{x}; \theta_o) \text{ for some } (\mathbf{y}, \mathbf{x}).$$

Conversely: if θ_o *not* identified, then some $\theta \neq \theta_o$ yields

$$f(\mathbf{y}|\mathbf{x}; \theta) = f(\mathbf{y}|\mathbf{x}; \theta_o) \text{ for all } (\mathbf{y}, \mathbf{x}).$$

data generating process would show the same densities
you would not be able to tell true θ and candidate θ 's apart

► I.e., θ and θ_o are **observationally equivalent**.

Example

Suppose that for some $(\alpha_o, \mu_o) \in \mathbb{R}^2$,

$$y_i = \alpha_o + \varepsilon_i, \quad \varepsilon_i \sim N(\mu_o, 1).$$

Q: Is (α_o, μ_o) identified?

no, they are conditionally dependent on each other
Would not be able separate them apart

A potential reason for non-convergence

Estimand

Identification implies (later): θ_o uniquely solves **population problem**

$$\max_{\theta \in \Theta} E [\ln f(\mathbf{y}_i | \mathbf{x}_i; \theta)] . \quad (\text{PP})$$

maximizing the expected log density

Equivalently, θ_o solves

$$\min_{\theta \in \Theta} E [-\ln f(\mathbf{y}_i | \mathbf{x}_i; \theta)] .$$

minimizing the negative log density

Taking $q(\mathbf{w}, \theta) = -\ln f(\mathbf{y} | \mathbf{x}; \theta) \Rightarrow \theta_o$ **M-estimand**.

Estimation

Analogy principle suggests **sample problem**

$$\max_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^N \ell_i(\boldsymbol{\theta}), \quad (\text{SP})$$

with **(conditional) likelihood contribution**

$$\ell_i(\boldsymbol{\theta}) := \ln f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}).$$

Maximum likelihood estimator (MLE): Any solution $\hat{\boldsymbol{\theta}}$ to SP.

► Every MLE an M-estimator!

Example: Probit

Example: Probit

Binary outcome y_i , i.e. $\mathcal{Y} = \{0, 1\}$,

$$p_o(y|\mathbf{x}) = p_o(1|\mathbf{x})^y [1 - p_o(1|\mathbf{x})]^{1-y}, \quad y \in \{0, 1\}.$$

Probit model $\{\Phi(\mathbf{x}\boldsymbol{\theta}); \boldsymbol{\theta} \in \Theta\}$ for $p_o(1|\mathbf{x})$.

► $\Phi: \mathbb{R} \rightarrow (0, 1)$: standard normal CDF.

► $\Theta \subseteq \mathbb{R}^P$.

Correctly specified if for some $\boldsymbol{\theta}_o \in \Theta$,

$$P(y_i = 1 | \mathbf{x}_i = \mathbf{x}) = \Phi(\mathbf{x}\boldsymbol{\theta}_o) \text{ for all } \mathbf{x} \in \mathcal{X}.$$

Example: Probit

Model implies density of $y_i|\mathbf{x}_i$,

$$f(y|\mathbf{x};\boldsymbol{\theta}) = \Phi(\mathbf{x}\boldsymbol{\theta})^y [1 - \Phi(\mathbf{x}\boldsymbol{\theta})]^{1-y}, \quad y \in \{0, 1\}.$$

Probit log-likelihood contribution

$$\ell_i(\boldsymbol{\theta}) = y_i \ln \Phi(\mathbf{x}_i\boldsymbol{\theta}) + (1 - y_i) \ln [1 - \Phi(\mathbf{x}_i\boldsymbol{\theta})].$$

Probit estimator solves

$$\max_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^N \{y_i \ln \Phi(\mathbf{x}_i\boldsymbol{\theta}) + (1 - y_i) \ln [1 - \Phi(\mathbf{x}_i\boldsymbol{\theta})]\}.$$

Identification and Solution Uniqueness

Identification and Uniqueness

Claim: θ_o identified \Rightarrow uniquely solves PP,

$$\theta_o = \operatorname{argmax}_{\theta \in \Theta} E [\ln f(\mathbf{y}_i | \mathbf{x}_i; \theta)] .$$

Will invoke **Jensen's inequality:** g concave + Z random

$$\Rightarrow E[g(Z)] \leq g(E[Z]) .$$

Inequality *strict* provided g *strictly* concave + Z *nonconstant*.

Identification and Uniqueness

- ▶ Fix $\theta \in \Theta$.
- ▶ $g := \ln(\cdot)$
- ▶ $Z := f(\mathbf{y}_i | \mathbf{x}_i; \theta) / f(\mathbf{y}_i | \mathbf{x}_i; \theta_o)$
- ▶ Cond'n on $\mathbf{x}_i \Rightarrow [\text{FILL IN}]$

Identification and Uniqueness

- ▶ Correct specification + legitimate density \Rightarrow

$$E \left[\frac{f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta})}{f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}_o)} \middle| \mathbf{x}_i \right] = \int_{\mathcal{Y}} \frac{f(\mathbf{y} | \mathbf{x}_i; \boldsymbol{\theta})}{f(\mathbf{y} | \mathbf{x}_i; \boldsymbol{\theta}_o)} p_o(\mathbf{y} | \mathbf{x}_i) \nu(d\mathbf{y})$$
$$=$$

- ▶ Hence

$$E \left[\ln \left(\frac{f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta})}{f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}_o)} \right) \middle| \mathbf{x}_i \right] \leq$$

- ▶ Rearranging,

$$E [\ln f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}_o) | \mathbf{x}_i] \geq E [\ln f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i].$$

Identification and Uniqueness

- ▶ Have shown:

$$E [\ln f (\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}_o) | \mathbf{x}_i] \geq E [\ln f (\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i] .$$

- ▶ Taking expectations,

$$E [\ln f (\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}_o)] \geq E [\ln f (\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta})] .$$

- ▶ $\boldsymbol{\theta} \in \Theta$ arbitrary $\Rightarrow \boldsymbol{\theta}_o$ solves PP.

Identification and Uniqueness

- ▶ Have shown: θ_o solves PP,

$$E [\ln f (\mathbf{y}_i|\mathbf{x}_i; \theta_o)] \geq E [\ln f (\mathbf{y}_i|\mathbf{x}_i; \theta)] \text{ for all } \theta \in \Theta.$$

- ▶ θ_o identified $\Rightarrow Z = f (\mathbf{y}_i|\mathbf{x}_i; \theta) / f (\mathbf{y}_i|\mathbf{x}_i; \theta_o)$ *nonconstant*.

- ▶ $\ln (\cdot)$ *strictly* concave, so Jensen \Rightarrow

$$E [\ln f (\mathbf{y}_i|\mathbf{x}_i; \theta_o)] > E [\ln f (\mathbf{y}_i|\mathbf{x}_i; \theta)] \text{ for all } \theta \in \Theta \setminus \{\theta_o\}.$$

- ▶ Hence, **identification implies unique maximizer**.

Asymptotic Properties

Asymptotic Properties of MLE

- ▶ Recall: Every MLE an M-estimator,

$$\hat{\boldsymbol{\theta}} \in \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^N -\ln f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}).$$

- ▶ May appeal to general results:
 - ▶ Consistency (W. Thm. 12.2)
 - ▶ Asymptotic normality (W. Thm. 12.3).
- ▶ Will verify relevant conditions.

Consistency

Recap: Consistency of M-Estimators

Theorem

If

1. $\Theta \subseteq \mathbb{R}^P$ compact (i.e. closed + bounded),
2. $q(\mathbf{w}, \cdot)$ continuous (in $\boldsymbol{\theta}$),
3. $\boldsymbol{\theta}_o$ uniquely minimizes $\boldsymbol{\theta} \mapsto E[q(\mathbf{w}_i, \boldsymbol{\theta})]$ (“identification”),

(+ technical conditions), then

1. Minimizer $\hat{\boldsymbol{\theta}}$ of $N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \cdot)$ exists,
2. $\hat{\boldsymbol{\theta}}$ consistent for $\boldsymbol{\theta}_o$.

Consistency of ML-Estimators

Q: Conditions verified?

1. Θ compact? Assumed...
2. $q(\mathbf{w}, \cdot) = -\ln f(\mathbf{y}|\mathbf{x}; \cdot)$ continuous?
 - ▶ Assume $f(\mathbf{y}|\mathbf{x}; \cdot)$ cont's.
 - ▶ **Probit:** Φ is cont's.
3. Unique PP solution?
 - ▶ Follows from θ_o identified.

$$\theta \neq \theta_o \Rightarrow f(\mathbf{y}|\mathbf{x}; \theta) \neq f(\mathbf{y}|\mathbf{x}; \theta_o) \text{ some } (\mathbf{y}, \mathbf{x}).$$

Compact + LL cont' + ML identification \Rightarrow MLE consistency.

Asymptotic Normality

Recap: Asymptotic Normality of **M**-Estimators

Theorem

Provided

- ▶ θ_o unique min'r + interior to Θ compact,
- ▶ $q(\mathbf{w}, \cdot)$ cont' + twice cont' diff' on $\text{int } \Theta$,
- ▶ $E[\mathbf{s}(\mathbf{w}_i, \theta_o)] = \mathbf{0}$, and $E[\mathbf{H}(\mathbf{w}_i, \theta_o)]$ positive definite,
- ▶ (+ technical),

$$\begin{aligned}\sqrt{N}(\hat{\theta} - \theta_o) &\xrightarrow{d} N(\mathbf{0}, \mathbf{A}_o^{-1} \mathbf{B}_o \mathbf{A}_o^{-1}), \\ \mathbf{A}_o &:= E[\mathbf{H}(\mathbf{w}_i, \theta_o)], \\ \mathbf{B}_o &:= E[\mathbf{s}(\mathbf{w}_i, \theta_o) \mathbf{s}(\mathbf{w}_i, \theta_o)'] .\end{aligned}$$

Oh, that Pesky Minus...

- ▶ Thm. 12.3 designed for *minimization*.
- ▶ Turn max'n into min'n:

$$q(\mathbf{w}_i, \boldsymbol{\theta}) = -\ln f(\mathbf{y}_i | \mathbf{x}_i; \boldsymbol{\theta}) = -\ell_i(\boldsymbol{\theta}).$$

- ▶ Hence above score (\mathbf{s})/Hessian (\mathbf{H}) of $-\ln f$ (wrt. $\boldsymbol{\theta}$).
- ▶ *In what follows,*

$$\mathbf{s}_i(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta})', \quad (\text{no minus})$$

$$\mathbf{H}_i(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \mathbf{s}_i(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}), \quad (\text{no minus})$$

$$\Rightarrow \mathbf{A}_o = -E[\mathbf{H}_i(\boldsymbol{\theta}_o)].$$

Information Matrix Equalities, I

- ▶ Additionally assuming...
 - ▶ θ_o interior to Θ ,
 - ▶ $\ln f(\mathbf{y}|\mathbf{x}; \cdot)$ twice cont' diff' on int Θ ,
 - ▶ (+ technical)
- ▶ May now apply Thm. 12.3.
- ▶ But further structure available...

Information Matrix Equalities, II

this is due to interchanging expectations (E) with derivatives
since expectations translate to integrals(?)

- ▶ Under quite mild conditions,

$$-E[\mathbf{H}_i(\boldsymbol{\theta}_o) | \mathbf{x}_i] = E[\mathbf{s}_i(\boldsymbol{\theta}_o) \mathbf{s}_i(\boldsymbol{\theta}_o)' | \mathbf{x}_i], \quad (\text{CIME})$$

$$\Rightarrow -E[\mathbf{H}_i(\boldsymbol{\theta}_o)] = E[\mathbf{s}_i(\boldsymbol{\theta}_o) \mathbf{s}_i(\boldsymbol{\theta}_o)']. \quad (\text{UIME})$$

- ▶ (Un)Conditional Information Matrix Equality.

- ▶ Implies $\mathbf{A}_o = \mathbf{B}_o$.

- ▶ Asymptotic variance simplifies.

Asymptotic Normality of ML-Estimators

Theorem

Provided

- ▶ θ_o identified + interior to Θ compact,
- ▶ $\ln f(\mathbf{y}|\mathbf{x}; \cdot)$ cont' + twice cont' diff' on $\text{int } \Theta$,
- ▶ + technical (including CIME justification),

$$\begin{aligned}\sqrt{N}(\hat{\theta} - \theta_o) &\xrightarrow{d} N(\mathbf{0}, \mathbf{A}_o^{-1}), \\ \mathbf{A}_o &:= -E[\mathbf{H}_i(\theta_o)].\end{aligned}$$

Hence $\text{Avar}(\hat{\theta}) = \mathbf{A}_o^{-1}/N$.

Asymptotic Variance Estimation

Asymptotic Variance Estimators

Three candidates for $\hat{\mathbf{A}}$:

$$= -\frac{1}{N} \sum_{i=1}^N \mathbf{H}_i(\hat{\boldsymbol{\theta}}), \quad (\text{least structure})$$

$$\text{or} \quad = \frac{1}{N} \sum_{i=1}^N \mathbf{s}_i(\hat{\boldsymbol{\theta}}) \mathbf{s}_i(\hat{\boldsymbol{\theta}})', \quad (\text{per UIME})$$

$$\text{or} \quad = \frac{1}{N} \sum_{i=1}^N \mathbf{A}(\mathbf{x}_i, \hat{\boldsymbol{\theta}}),$$

where $\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o) := -E[\mathbf{H}_i(\boldsymbol{\theta}_o) | \mathbf{x}_i]$.

Each $\hat{\mathbf{A}} \rightarrow_p \mathbf{A}_o (= \mathbf{B}_o)$ under mild (add'l) cond's.

Avar Estimation: Discussion

$$\widehat{\text{Avar}}(\hat{\boldsymbol{\theta}}) = \underbrace{\left(-\sum_{i=1}^N \mathbf{H}_i(\hat{\boldsymbol{\theta}})\right)^{-1}}_{(1)}, \underbrace{\left(\sum_{i=1}^N \mathbf{s}_i(\hat{\boldsymbol{\theta}})\mathbf{s}_i(\hat{\boldsymbol{\theta}})'\right)^{-1}}_{(2)}, \text{ or } \underbrace{\left(\sum_{i=1}^N \mathbf{A}(\mathbf{x}_i, \hat{\boldsymbol{\theta}})\right)^{-1}}_{(3)} \quad ?$$

Pros/Cons:

1. Always available, 2nd-order diff', p.s.d.
2. Easy to compute, 1st-order diff', p.s.d.
3. Harder to derive, often p.d. + good in small sample.

Example: Probit Avar Estimation

Probit Avar Estimation

Recall: Cond'l probit density

$$f(y|\mathbf{x};\boldsymbol{\theta}) = \Phi(\mathbf{x}\boldsymbol{\theta})^y [1 - \Phi(\mathbf{x}\boldsymbol{\theta})]^{1-y}, \quad y \in \{0, 1\}.$$

Cond'l probit LL:

$$\ell_i(\boldsymbol{\theta}) = y_i \ln \Phi(\mathbf{x}_i\boldsymbol{\theta}) + (1 - y_i) \ln [1 - \Phi(\mathbf{x}_i\boldsymbol{\theta})].$$

We'll use option (3) (“cond'l Hessian”):

1. Derive score, $\mathbf{s}_i(\boldsymbol{\theta})$.
2. Derive $\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o) = -E[\mathbf{H}_i(\boldsymbol{\theta}_o)|\mathbf{x}_i]$.
3. Sum + plug in $\hat{\boldsymbol{\theta}}$ + invert.

Probit Avar Estimation

Step 1: Derive score.

Chain rule + gather \Rightarrow

$$\mathbf{s}_i(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}) = \frac{[y_i - \Phi(\mathbf{x}_i \boldsymbol{\theta})] \varphi(\mathbf{x}_i \boldsymbol{\theta})}{\Phi(\mathbf{x}_i \boldsymbol{\theta}) [1 - \Phi(\mathbf{x}_i \boldsymbol{\theta})]} \mathbf{x}_i'.$$

Step 2: Derive $\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o)$.

$$\begin{aligned} \mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o) &= -E[\mathbf{H}_i(\boldsymbol{\theta}_o) | \mathbf{x}_i] \\ &= E[\mathbf{s}_i(\boldsymbol{\theta}_o) \mathbf{s}_i(\boldsymbol{\theta}_o)' | \mathbf{x}_i]. \quad (\text{CIME}) \\ &= E\left\{ \frac{[y_i - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)]^2 \varphi(\mathbf{x}_i \boldsymbol{\theta}_o)^2}{\Phi(\mathbf{x}_i \boldsymbol{\theta}_o)^2 [1 - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)]^2} \mathbf{x}_i' \mathbf{x}_i \middle| \mathbf{x}_i \right\}. \end{aligned}$$

outer product product are like squaring

Probit Avar Estimation

Step 2: Derive $\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o)$ (ctnd).

$$\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o) = E \left\{ [y_i - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)]^2 \middle| \mathbf{x}_i \right\} \frac{\varphi(\mathbf{x}_i \boldsymbol{\theta}_o)^2}{\Phi(\mathbf{x}_i \boldsymbol{\theta}_o)^2 [1 - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)]^2} \mathbf{x}_i' \mathbf{x}_i.$$

$$y_i \text{ binary} + p_o(1|\mathbf{x}_i) = \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)$$

$$\Rightarrow E \left\{ [y_i - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)]^2 \middle| \mathbf{x}_i \right\} = \Phi(\mathbf{x}_i \boldsymbol{\theta}_o) [1 - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)].$$

Hence

$$\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o) = \frac{\varphi(\mathbf{x}_i \boldsymbol{\theta}_o)^2}{\Phi(\mathbf{x}_i \boldsymbol{\theta}_o) [1 - \Phi(\mathbf{x}_i \boldsymbol{\theta}_o)]} \mathbf{x}_i' \mathbf{x}_i$$

Probit Avar Estimation

Step 3: Sum + plug in $\hat{\theta}$ + invert.

\Rightarrow variance matrix estimator for probit:

$$\widehat{\text{Avar}}(\hat{\theta}) = \left(\sum_{i=1}^N \mathbf{A}(\mathbf{x}_i, \hat{\theta}) \right)^{-1},$$

$$\mathbf{A}(\mathbf{x}_i, \hat{\theta}) := \frac{\varphi(\mathbf{x}_i' \hat{\theta})^2}{\Phi(\mathbf{x}_i' \hat{\theta})[1 - \Phi(\mathbf{x}_i' \hat{\theta})]} \mathbf{x}_i' \mathbf{x}_i.$$

- ▶ Positive definite when invertible.
- ▶ Same result possible from 2nd-order diff'n. (Check!)