Econometrics I Lecture 9: Maximum Likelihood Estimation

Paul T. Scott NYU Stern

Fall 2021

Paul T. Scott NYU Stern

Econometrics I

Fall 2021 1 / 38

• What's a **likelihood**? It's basically the probability of the data conditional on a parameter value *θ*:

```
Pr (observed data|\theta),
```

but we think of this as a function of θ and telling us something about the plausibility of $\theta.$

- Computing likelihoods requires a model of what the probability of the data is.
- ⇒ In comparison to GMM estimation, Likelihood-based estimation requires relatively strong assumptions about the data generating process.

• Bayes's rule:

$$Pr\left(heta| extsf{data}
ight) = rac{Pr\left(extsf{data}| heta
ight)Pr\left(heta
ight)}{Pr\left(extsf{data}
ight)}$$

- Here, we directly make judgments about the (relative) probabilities of different parameter values.
- $Pr(data|\theta)$ is a **likelihood**
- Pr(data) is kind of irrelevant it's the same for any θ
- Pr (θ) is a prior. If we ignore it or if we assume it's the same for all θ, then we're just in the MLE world.

Likelihood Function

• Let $f(\cdot|\theta)$ represent the probability density of the data conditional on a parameter value θ . If data are independently and identically distributed, the **likelihood function** is

$$L(\boldsymbol{\theta}|\mathbf{y}) = f(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n|\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{y}_i|\boldsymbol{\theta})$$

where \mathbf{y}_i indicates individual observations (including both dependent and explanatory variables).

• We typically work with **log-likelihood function** because it's computationally simpler:

$$\ln L(\boldsymbol{\theta}|\mathbf{y}) = \sum_{i=1}^{n} \ln f(\mathbf{y}_{i}|\boldsymbol{\theta}).$$

• Maximum likelihood estimation entails estimating θ by maximizing the likelihood function:

$$\hat{oldsymbol{ heta}} = rg \min_{oldsymbol{ heta}} L\left(oldsymbol{ heta}| \mathbf{y}
ight) = rg \min_{oldsymbol{ heta}}$$
 In $L\left(oldsymbol{ heta}| \mathbf{y}
ight)$

• Since the natural log function is strictly increasing, maximizing the likelihood and maximizing log likelihood amount to the same thing.

• Recall that PDF of the normal distribution is

$$f_{\mathcal{N}}\left(arepsilonert\sigma
ight) = rac{1}{\sqrt{2\pi\sigma^2}}\exp\left(rac{-arepsilon^2}{2\sigma^2}
ight)$$

(for normal ε with zero mean and variance σ^2)

• Thus, log likelihood of an individual observation of ε_i is

$$\ln f_{\mathcal{N}}(\varepsilon_i | \sigma) = -\frac{1}{2} \left(\ln \sigma^2 + \ln 2\pi + \frac{\varepsilon_i^2}{\sigma^2} \right)$$

Likelihood for Linear Regression Model

 For linear model with ε_i mean-zero normal conditional on x_i, the likelihood of one observation is

$$L(\boldsymbol{\beta}, \sigma | y_i, \mathbf{x}_i) = f_{\mathcal{N}} \left(y_i - \mathbf{x}'_i \boldsymbol{\beta} | \sigma \right)$$

noting that this requires the distribution of ε_i to be mean-zero normal *conditional* on \mathbf{x}_i .

• Assuming the data are i.i.d across observations, the conditional likelihood of all the data is then

$$\ln L(\boldsymbol{\beta}, \sigma | \mathbf{y}, \mathbf{X}) = \sum_{i=1}^{n} \ln f_{\mathcal{N}} \left(y_{i} - \mathbf{x}_{i}^{\prime} \boldsymbol{\beta} | \sigma \right)$$
$$= -\frac{1}{2} \sum_{i=1}^{n} \left(\ln \sigma^{2} + \ln 2\pi + \frac{\left(y_{i} - \mathbf{x}_{i}^{\prime} \boldsymbol{\beta} \right)^{2}}{\sigma^{2}} \right)$$

MLE for Linear Model I

• Linear model log-likelihood:

$$\ln L(\boldsymbol{\beta}, \sigma | \mathbf{y}, \mathbf{X}) = -\frac{1}{2} \sum_{i=1}^{n} \left(\ln \sigma^2 + \ln 2\pi + \frac{(y_i - \mathbf{x}'_i \boldsymbol{\beta})^2}{\sigma^2} \right)$$

• Focus on the term that involves β :

$$\frac{-1}{2\sigma^2}\sum_{i=1}^n \left(y_i - \mathbf{x}_i'\beta\right)^2$$

NB: maximizing the likelihood with respect to β is equivalent to least squares

MLE for Linear Model II

- MLE estimate of β is the same as OLS.
- MLE estimate of σ^2 comes from setting $\frac{d}{d\sigma} \ln L\left(\hat{\boldsymbol{\beta}}, \sigma | \mathbf{y}, \mathbf{X}\right) = 0$:

$$\hat{\sigma}_{MLE}^2 = n^{-1} \sum_{i=1}^n e_i^2$$

where $e_i = y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}$.

• Note that this is a bit different than the estimate of σ^2 we saw before:

$$s^{2} = (n - K)^{-1} \sum_{i=1}^{n} e_{i}^{2}$$

but the difference will be small in large samples. Recall: s^2 is a unbiased estimate of σ^2 , so this means that the ML estimate is biased, but very slightly in large samples.

Paul T. Scott NYU Stern

- An estimator is **asymptotically efficient** if its asymptotic covariance matrix is not larger than any other consistent estimator (i.e., standard errors are as small as any other estimator).
- It can be shown that (under regularity conditions), MLE is asymptotically efficient.
- Thus, MLE always performs well in large samples.

Estimating Standard Errors I

 The first way to estimate the asymptotic covariance matrix is to take second derivatives of the likelihood function:

$$\mathbf{\Gamma}^{-1} = \left(-\frac{\partial^2 \ln L\left(\hat{\boldsymbol{\theta}}\right)}{\partial \hat{\boldsymbol{\theta}} \partial \hat{\boldsymbol{\theta}}'}\right)^{-1}$$

• A second way is to compute the covariance of the first derivatives:

$$\mathbf{S}^{-1} = \left[\sum_{i=1}^{n} \hat{\mathbf{g}}_i \hat{\mathbf{g}}'_i\right]^{-1}$$

where

$$\hat{\mathbf{g}}_i = rac{\partial \ln f\left(\mathbf{x}_i, \hat{\boldsymbol{ heta}}\right)}{\partial \hat{\boldsymbol{ heta}}}.$$

• Either of the above is an asymptotically consistent estimator of $V(\hat{\theta}_{MLE})$. The latter is usually easier to compute.

Paul T. Scott NYU Stern

Econometrics I

MLE as GMM

• To maximize the likelihood function we set

$$n^{-1}\sum_{i=1}^{n}\hat{\mathbf{g}}_{i}=n^{-1}\sum_{i=1}^{n}\frac{\partial \ln f\left(\mathbf{x}_{i},\hat{\boldsymbol{\theta}}\right)}{\partial \hat{\boldsymbol{\theta}}}=0.$$

Thus, maximum likelihood is a GMM estimator based on moments

$$E\left[\frac{\partial \ln f\left(\mathbf{x}_{i}, \hat{\boldsymbol{\theta}}\right)}{\partial \hat{\boldsymbol{\theta}}}\right] = 0.$$

• The GMM estimator for the asymptotic covariance matrix has the form

$$\left(\mathbf{\Gamma} \mathbf{S}^{-1} \mathbf{\Gamma}
ight)^{-1}$$
 ,

but in the MLE context it can be shown that **S** and **Γ** are asymptotically equivalent, so they effectively cancel and we can use either S^{-1} or Γ^{-1} to estimate the variance.

Paul T. Scott NYU Stern

• Our starting point was that likelihoods were about the probability of the data conditional on a parameter value:

$$\ln L(\theta | \text{data}) = \sum_{i=1}^{n} \ln f(\text{data}_{i} | \theta).$$

- The above derivation was about ε_i, or the probability of y_i|x_i. But x_i might be a random variable, and it's also part of the data.
- Do we need to consider the randomness in x_i? In econometric models, typically we don't bother to explicitly model the randomness in explanatory variables.

Conditional Likelihood II

• Start with the full log likelihood function

$$\sum_{i=1}^n \ln p(y_i, \mathbf{x}_i | \alpha)$$

• We can decompose this using $Pr(y_i, \mathbf{x}_i) = Pr(y_i | \mathbf{x}_i) Pr(\mathbf{x}_i)$:

$$\sum_{i=1}^{n} \ln f(y_i | \mathbf{x}_i, \theta) + \sum_{i=1}^{n} \ln g(\mathbf{x}_i, \delta)$$

where θ is the subset of α that dictates the distribution of $y_i | \mathbf{x}_i$ and δ is the subset of α that dictates the distribution of \mathbf{x}_i .

 If we're only interested in θ, then as long as there are no restrictions between θ and δ, we can just focus on the first component of the likelihood function (i.e., the conditional likelihood function)

Conditional Likelihood III

• We can decompose this using $Pr(y_i, \mathbf{x}_i) = Pr(y_i | \mathbf{x}_i) Pr(\mathbf{x}_i)$:

$$\sum_{i=1}^{n} \ln f(y_i | \mathbf{x}_i, \theta) + \sum_{i=1}^{n} \ln g(\mathbf{x}_i, \delta)$$

where θ is the subset of α that dictates the distribution of $y_i | \mathbf{x}_i$ and δ is the subset of α that dictates the distribution of \mathbf{x}_i .

- Bottom line: you don't always have to specify and estimate a complete data generating process to do maximum likelihood estimation.
- Sometimes g and δ are of interest in any case (e.g., for counterfactual simulations).

- Note that the likelihood framework does not solve the endogeneity problem.
- The consistency of MLE relies on the model being correctly specified, and when ε_i and x_i are correlated, the mean of ε_i is generally non-zero conditional on x_i.
- Full information maximum likelihood (FIML) and limited information maximum likelihood (LIML) are the ML analog of IV estimators.
 - Because they require specifying a distribution for the error terms (typically normal) while 2SLS and GMM regression do not, ML-based IV estimators are not popular.
 - One advantage: LIML is more robust to weak instruments than 2SLS.

- Censored data is a common problem
 - Demand for a concert/sporting event with capacity constraints.
 - Meters often only measure outcomes within a bounded range (speedometers, thermometers, etc.)
 - ► A test is scored on a bounded range (200-800), and we're thinking of the test as marker for ability.

ML Application 1: Censored Regression Model II



How would we go about estimating this model?

Paul T. Scott NYU Stern

Econometrics I

- Suppose *v* is distributed with standard normal PDF, but only for values above a cutoff *a*.
- PDF will be

$$\frac{\phi\left(v\right)}{1-\Phi\left(a\right)}$$

where ϕ is the standard normal PDF and Φ is standard normal CDF.

• Note that we must divide by $1 - \Phi(a)$ to make the PDF integrate to 1.

Truncated Normal Moments I

Truncated Normal Properties

Suppose $v \sim \mathcal{N}(0, 1)$ has a normal distribution truncated with v > a. That is, v takes values in (a, ∞) and has PDF

$$\frac{\phi\left(v\right)}{1-\Phi\left(a\right)}$$

Then,

$$E[v] = \frac{\phi(a)}{1-\Phi(a)}$$

Var [v] = $\left(1 - \frac{\phi(a)}{1-\Phi(a)} \left(\frac{\phi(a)}{1-\Phi(a)} - a\right)\right)$

The ratio of a normal density to its CDF, $\frac{\phi(v)}{1-\Phi(a)}$, is known as the **inverse** Mills ratio.

• If original distribution is $v \sim \mathcal{N}(\mu, \sigma^2)$, truncated for v > a, we get similar results:

$$E[v] = \mu + \sigma \frac{\phi(\alpha)}{1 - \Phi(\alpha)}$$

$$Var[v] = \sigma^2 \left(1 - \frac{\phi(\alpha)}{1 - \Phi(\alpha)} \left(\frac{\phi(\alpha)}{1 - \Phi(\alpha)} - \alpha\right)\right)$$

where $\alpha = \frac{a-\mu}{\sigma}$.

• If truncation is for v < a, then we replace $\frac{\phi(\alpha)}{1-\Phi(\alpha)}$ with $-\frac{\phi(\alpha)}{\Phi(\alpha)}$

• Suppose $v^* \sim \mathcal{N}(\mu, \sigma^2)$. Consider

$$v = egin{cases} v^* & ext{if } v^* > a \ a & ext{if } v^* \leq a \end{cases}$$

 Note: v will have the normal PDF above the cutoff a, and there will be a point mass at v = a.

• $Pr(v = a) = \Phi\left(\frac{a-\mu}{\sigma}\right)$ where Φ is the standard normal CDF.

Censored Normal Mean

• Censored Normal will have mean

$$E(v) = E(v|v = a) Pr(v = a) + E(v|v > a) Pr(v > a)$$
$$= a\Phi + E(v|v > a)(1 - \Phi)$$
$$= a\Phi + (\mu + \sigma\lambda)(1 - \Phi)$$

where
$$\lambda = \frac{\phi(\alpha)}{1-\Phi(\alpha)}$$
, $\Phi = \Phi(\alpha)$, $\alpha = \frac{a-\mu}{\sigma}$

• We can similarly derive the variance from the truncated normal variance

$$Var(v) = \sigma^{2}(1-\Phi)\left[(1-\delta) + (\alpha-\lambda)^{2}\Phi\right]$$

where $\delta = \lambda^2 - \lambda \alpha$.

• Let's now return to censored regression framework:

$$\begin{array}{rcl} y_i^* = & \mathbf{x}_i \boldsymbol{\beta} + \varepsilon_i \\ y_i = & 0 & \quad \text{if } y_i^* \leq 0 \\ y_i = & y_i^* & \quad \text{if } y_i^* > 0 \end{array}$$

- What do you expect to happen if we estimate with OLS?
- What if we drop the observations with $y_i = 0$?

Censored Regression: Bias in OLS

$$\begin{array}{rcl} y_i^* = & \mathbf{x}_i \boldsymbol{\beta} + \varepsilon_i \\ y_i = & 0 & \quad \text{if } y_i^* \leq 0 \\ y_i = & y_i^* & \quad \text{if } y_i^* > 0 \end{array}$$

 If we run OLS on the censored regression data, we are essentially estimating

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + \varepsilon_i^*$$

where $\varepsilon_i^* = \varepsilon_i + y_i - y_i^*$

 Is ε^{*}_i uncorrelated with x as we would need for OLS to deliver an unbiased estimate?

Censored Regression: ML Estimation (Tobit)

• Log likelihood equation:

$$\ln L = \sum_{y_i > 0} -\frac{1}{2} \left[\ln \left(2\pi \right) + \ln \sigma^2 + \frac{\left(y_i - \mathbf{x}'_i \boldsymbol{\beta} \right)^2}{\sigma^2} \right] + \sum_{y_i = 0} \ln \left(1 - \Phi \left(\frac{\mathbf{x}'_i \boldsymbol{\beta}}{\sigma} \right) \right)$$

- Maximum likelihood here will give consistent (and asymptotically efficient) estimates of all parameters.
- This is known as a **tobit** regression.
- These mathematical tools are also what's behind the **Heckman** selection correction to deal with *sample selection bias*.

Application 2: Finite Mixture Models

- x observed variables
- ζ unobserved variables assumed to have finite support, Z
- θ parameters of interest
- $p(x_i, \zeta_i | \theta)$ complete data likelihood for *i*th observation
- $p(x_i|\theta)$ incomplete data likelihood for *i*th observation:

$$p(x_i|\theta) = \sum_{z \in Z} p(x_i, z|\theta)$$

• $q_{iz}\left(heta
ight)$ - expectation of incomplete data

$$q_{iz}\left(\theta\right) = \Pr\left(\zeta_{i} = z | x_{i}, \theta\right)$$

Example 1: Mixture of Normals

• $Pr(z_i = 1) = \alpha_1$



Example 2: collusion (Porter, 1983)

• Rob Porter (1983), "A Study of Cartel Stability: The Joint Executive Committee, 1880-1886"

$$\begin{aligned} \ln Q_t &= \alpha_0 + \alpha_1 \ln P_t + \alpha_2 D_t + U_{1t} \\ \ln P_t &= \beta_0 + \beta_1 \ln Q_t + \beta_2 S_t + \beta_3 I_t + U_{2t} \end{aligned}$$

where

- D_t: demand shifters
- St: supply shifters
- ▶ $I_t \in \{0, 1\}$ indicating whether the cartel was in a price war or not
- In previous notation,
 - $\flat x_t = (Q_t, P_t, D_t, S_t)$
 - $\blacktriangleright z_t = I_t$
 - θ = (α, β)
 - ▶ to deal with simultaneity, likelihood function $p(x_i, \zeta_i | \theta)$ is FIML

The *incomplete data log-likelihood function* or *unconditional log-likelihood function* for a mixture model involves a sum within an expectation, which makes it very hard to maximize with standard optimization algorithms:

$$\ln L(x|\theta) = \sum_{i} \ln \left(\sum_{z} p(x_{i}, z|\theta) \right).$$

The EM algorithm is based on the (expected) *complete data log-likelihood function*:

$$Q(x, q|\theta) = \sum_{i} \sum_{z} q_{iz} \ln (p(x_i, z|\theta)).$$

Note that Q would simply be the log-likelihood function if ζ were observed.

EM Algorithm overview

- The EM algorithm starts with some initial guess for $\theta^{(0)}$
- In the E-step, we calculate expectations of the *q*'s conditional on the parameter values:

$$q_{iz}^{(m)} = Pr\left(\zeta_i = z|\theta^{(m-1)}\right).$$

• In the M-step, we maximize the value of the complete data likelihood function:

$$\theta^{(m)} = \max_{\theta} Q\left(x, q^{(m)}|\theta\right).$$

• The EM Algorithm iteratively applies E and M steps until $\theta(m)$ converges.

EM Algorithm overview

- The E and M steps are often easy computationally (in contrast to maximization of incomplete data likelihood function).
- Each EM iteration increases $Q\left(x, q^{(m)}|\theta^{(m)}\right)$, and $L(x|\theta) \ge Q\left(x, q^{(m)}|\theta^{(m)}\right)$.
- Thus, iterating on the E and M steps will monotonically increase $\ln L\left(x|\theta^{(m)}\right)$, and $\theta^{(m)}$ will generally converge to a local maximum of $\ln L\left(x|\theta\right)$.
- ⇒ EM Algorithm transforms a hard optimization problem into a series of easy optimization problems

Estimation of Mixture of Normals I

•
$$\theta = (\mu_1, \mu_2, \sigma, \alpha_1)$$

• If $z_i = 1$, then $x_i \sim N(\mu_1, \sigma)$
• If $z_i = 2$, then $x_i \sim N(\mu_2, \sigma)$
• $Pr(z_i = 1) = \alpha_1$

In the E step, we just apply Bayes's Theorem to find q's

$$q_{i1}^{(m)} = \Pr\left(z_i = 1 | x_i, \theta^{(m)}\right) = \alpha_1^{(m)} f\left(x_i | \mu_1^{(m)}, \sigma^{(m)}\right) \\ \frac{\alpha_1^{(m)} f\left(x_i | \mu_1^{(m)}, \sigma^{(m)}\right) + \left(1 - \alpha_1^{(m)}\right) f\left(x_i | \mu_2^{(m)}, \sigma^{(m)}\right)}{\alpha_1^{(m)} f\left(x_i | \mu_1^{(m)}, \sigma^{(m)}\right) + \left(1 - \alpha_1^{(m)}\right) f\left(x_i | \mu_2^{(m)}, \sigma^{(m)}\right)}$$

where $f(x|\mu, \sigma)$ is the density at x of the normal distribution with mean μ and standard deviation σ^2 .

Estimation of Mixture of Normals II

• In the M step, maximizing the complete data likelihood function amounts to taking weighted means:

$$\mu_z^{(m)} = \sum_i q_{iz}^{(m)} x_i$$

$$\sigma^{(m)} = \sqrt{\frac{\sum_{z} \sum_{i} q_{iz}^{(m)} (x_{i} - \mu_{z})^{2}}{\sum_{z} \sum_{i} q_{iz}^{(m)}}}$$
$$\alpha_{z}^{(m)} = N^{-1} \sum_{i} q_{iz}^{(m)}$$

Estimation of example 1: mixture of normals

- Note: in a mixture model with covariates that enter linearly, the M step involves weighted OLS instead of a weighted mean
- Bottom line: E and M step are both easy computationally, so iterating on them goes quickly.
- In general, the EM algorithm can stop at local maxima, so some care is needed to ensure a global optimum is attained (e.g., multiple starting points).
- Alternative: MCMC estimation.

Model Selection: Likelihood Ratio

- When comparing nested models, the likelihood ratio test is simple and powerful
- Let heta be a vector of parameters to be estimated
 - $\hat{\theta}_U$ is the ML estimate for the full model
 - $\hat{\theta}_R$ is the ML estimate for a restricted model (e.g., with a couple elements fixed to zero)
- Likelihood ratio:

$$\lambda = rac{L\left(\hat{oldsymbol{ heta}}_R|\mathsf{data}
ight)}{L\left(\hat{oldsymbol{ heta}}_U|\mathsf{data}
ight)},$$

which will always be less than one.

- Null hypothesis H_0 : the restricted model is correct.
- Given regularity conditions and H_0 , then asymptotically asymptotic distribution of

–2 ln
$$\lambda \sim \mathcal{X}_R^2$$
,

where \mathcal{X}^2_R is chi-squared distribution with degrees of freedom equal to number of restrictions.

• Note similarly to testing restrictions in linear models, but no need for linearity and computationally simpler than F test.

Model Selection: Information Criteria

- Just as R^2 always increases as we add parameters, so does the likelihood.
- When comparing models with different numbers of parameters, we should penalize more complex models. Intuitively, evaluating models based on likelihood without a penalty will lead to over fitting the data.
- Two popular criteria for selecting models that reward parsimony:

Akaike information criterion $= -2 \ln L(\theta | \mathbf{y}) + 2K$ Bayes information criterion $= -2 \ln L(\theta | \mathbf{y}) + K \ln n$

- To compare two or more models using the AIC (BIC), compute each model's AIC (BIC) score, and select the model with the lowest score (highest penalized likelihood).
- Note: these can be used to compare non-nested models as well as nested models.