

Nonlinear estimation and M-estimation

- Nonlinear estimation describes any problem in which estimators cannot be obtained in closed form
- M-estimators (M for maximisation) is a vary wide class of estimators. It includes
 - maximum likelihood
 - nonlinear least squares
 - least absolute deviation (*LAD*)
 - quasi-maximum likelihood
 - and many others

Nonlinear Least Squares (*NLS*)

- Nonlinear regression function. For some $\theta_o \in \Theta$

$$E(y | \mathbf{x}) = m(\mathbf{x}, \theta_o)$$

where $m(\mathbf{x}, \theta)$ is a known function and θ_o is a finite parameter vector

- Function $m(\mathbf{x}, \theta_o)$ is describing the nonlinear dependence between the conditional expected values of the dependent variable and explanatory variables

Example 1. *Sometimes the linear model is not well suited to the data we have. If $y > 0$ (e.g. incomes) or $y \in [0, 1]$ (e.g. percentage of unemployed) the*

linear model is problematic as for $E(y|\mathbf{x}) = \mathbf{x}\beta$ the values of $y \in (-\infty, \infty)$ for $\mathbf{x} \in R^K$.

For $y > 0$ exponential regression function can be used $m(\mathbf{x}, \boldsymbol{\theta}) = \exp(\mathbf{x}\boldsymbol{\theta})$. For $y \in [0, 1]$ the logistic model $m(\mathbf{x}, \boldsymbol{\theta}) = \frac{\exp(\mathbf{x}\boldsymbol{\theta})}{1 + \exp(\mathbf{x}\boldsymbol{\theta})}$ is often used.

- We will denote true parameter vector as $\boldsymbol{\theta}_o$
- The model we consider is nonlinear regression model in error form

$$y = m(\mathbf{x}, \boldsymbol{\theta}_o) + u$$

and $E(u|\mathbf{x}) = 0$

- Let $\mathbf{w}_i \equiv (\mathbf{x}_i, y_i)$ be the vector of observation drawn from random sample

- Our estimation is based on the fact that θ_o minimizes the squared error in the population

$$\min_{\theta \in \Theta} \left\{ E [y - m(\mathbf{x}, \theta_o)]^2 \right\}$$

- The estimation is based on the analogy principle we calculate estimator of θ_o by minimizing the squared error in the sample

$$\min_{\theta \in \Theta} N^{-1} \sum_{i=1}^N [y_i - m(\mathbf{x}, \theta)]^2$$

General case

- The M-estimator $\hat{\theta}$ solves the

$$\min_{\theta \in \Theta} N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \theta)$$

for some function $q(\mathbf{w}_i, \theta)$

- The parameter vector θ_o is assumed to solve uniquely the population problem (is identified)

$$\min_{\theta \in \Theta} E[q(\mathbf{w}, \theta)]$$

- Because of the law of large numbers

$$N^{-1} \sum_{i=1}^N q(\mathbf{w}_i, \boldsymbol{\theta}) \xrightarrow{p} E[q(\mathbf{w}_i, \boldsymbol{\theta})]$$

and as $\hat{\boldsymbol{\theta}}$ is minimizing the term on left and $\boldsymbol{\theta}_o$ the term on the right it seems plausible that $\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_o$

Identification

- The identification in the context of M-estimation means that θ_o is a unique solution to the population problem
- For *NLS* it means that

$$E \{ [m(\mathbf{x}, \theta_o) - m(\mathbf{x}, \theta)] \} > 0, \text{ for all } \theta \in \Theta, \theta \neq \theta_o$$

Example 2. *In regression function $m(\mathbf{x}, \theta) = \theta_1 + \theta_2 x_2 + \theta_3 x_3^{\theta_4}$ this condition fails if $\theta_3 = 0$ (poorly identified model)*

Asymptotic normality of M-estimators

- The estimator $\hat{\theta}$ minimizes the expected value of $q(\mathbf{w}, \theta)$
- Then if $q(\mathbf{w}, \theta)$ is continuously differentiable w.r.t. θ than the first order derivative of $q(\mathbf{w}, \theta)$ should be equal to zero at $\hat{\theta}$, as $\hat{\theta}$ solves the maximisation problem
- Denoted this first order derivative as $s_i = s(\mathbf{w}_i, \theta) = \frac{\partial q(\mathbf{w}_i, \theta)}{\partial \theta}$. Then

$$\sum_{i=1}^N s(\mathbf{w}_i, \hat{\theta}) = 0$$

- Function $s(\mathbf{w}, \theta)$ is called score of the objective function

- Define $\ddot{H}_i = \frac{\partial^2 q(\mathbf{w}_i, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}$ as Hessian of the objective function for observation i
- Assume that \ddot{H}_i fulfills the conditions of the law of large numbers and so $\frac{1}{N} \sum_{i=1}^N \ddot{H}_i \xrightarrow{p} E[\mathbf{H}(\mathbf{w}, \boldsymbol{\theta}_o)] = \mathbf{A}_o$ (and assume that \mathbf{A}_o is nonsingular)
- Assume that $\mathbf{s}(\mathbf{w}_i, \boldsymbol{\theta}_o)$ satisfies the conditions of central limit theorem as a *i.i.d.* random variable with mean zero and then $N^{-\frac{1}{2}} \sum_{i=1}^N \mathbf{s}(\mathbf{w}_i, \boldsymbol{\theta}_o) \xrightarrow{D} N(0, \mathbf{B}_o)$ where $\mathbf{B}_o = E[\mathbf{s}(\mathbf{w}, \boldsymbol{\theta}_o) \mathbf{s}(\mathbf{w}, \boldsymbol{\theta}_o)']$
- Under these conditions

$$\sqrt{N} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_o) \xrightarrow{D} N(0, \mathbf{A}_o^{-1} \mathbf{B}_o \mathbf{A}_o^{-1})$$

- The asymptotic variance can then be estimated as

$$Avar(\hat{\boldsymbol{\theta}}) = \frac{1}{N} \mathbf{A}_o^{-1} \mathbf{B}_o \mathbf{A}_o^{-1}$$

Estimating the variance covariance matrix

- In practice it difficult to calculate directly $B_o = E [s (w, \theta_o) s (w, \theta_o)']$ and $A_o = E [H (w, \theta_o)]$.
- We can however estimate $H (w, \theta_o)$ using the Law of Large Numbers:

$$\hat{A} = N^{-1} \sum_{i=1}^N H (w_i, \hat{\theta}) = N^{-1} \sum_{i=1}^N \hat{H}_i \xrightarrow{p} A_o$$

- Drawbacks of this estimator:
 - need to derive the second order derivatives
 - this estimator is not always positive definite

- In most economic application the inference is conditional on \mathbf{x}_i then we use the conditional variance matrix

$$\mathbf{A}_i = \mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o) = \mathbb{E}[\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}) | \mathbf{x}_i]$$

- But by the law of iterated expectations $\mathbb{E}[\mathbf{A}(\mathbf{x}_i, \boldsymbol{\theta}_o)] = \mathbb{E}[\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta})] = \mathbf{A}_o$

- Again

$$\widehat{\mathbf{A}} = N^{-1} \sum_{i=1}^N \mathbf{A}(\mathbf{x}_i, \widehat{\boldsymbol{\theta}}) = N^{-1} \sum_{i=1}^N \widehat{\mathbf{A}}_i \xrightarrow{p} \mathbf{A}_o$$

- This estimator is useful if $\mathbb{E}[\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}) | \mathbf{x}_i]$ can be obtained. In most cases it is possible to prove that this estimator is positive definite.

- We obtain the consistent estimator of B_o in the similar way

$$\widehat{B} = N^{-1} \sum_{i=1}^N s(w_i, \widehat{\theta}) s(w_i, \widehat{\theta})' = N^{-1} \sum_{i=1}^N \widehat{s}_i \widehat{s}_i' \xrightarrow{p} B_o$$

- Then the consistent estimator of the variance covariance matrix of $\sqrt{N}(\widehat{\theta} - \theta_o)$ has the form

$$\widehat{V} = \widehat{A}^{-1} \widehat{B} \widehat{A}^{-1}$$

- And the estimate of the variance of $\widehat{\theta}$

$$Av\hat{a}r(\widehat{\theta}) = \widehat{A}^{-1} \widehat{B} \widehat{A}^{-1} / N$$

- In many contexts it is possible to simplify this formula under additional assumption

$$E [s (w, \theta) s (w, \theta)'] = \sigma_o^2 E [H (w, \theta_o)]$$

or simply $B_o = \sigma^2 A_o$

- This assumption is called generalized information matrix equality (GIME)
- Under this assumption the variance estimator has the form

$$\hat{V} = \hat{\sigma}^2 \hat{A}^{-1}$$

Hypothesis Testing

- We have three test for testing nonlinear hypothesis $H_0 : c(\theta_o) = 0$
 - Wald
 - Score (Lagrange Multiplier test)
 - Tests based on the change of the objective function (Quasi Likelihood Ratio test)
- $C(\theta_0) = \frac{\partial c(\theta_o)}{\partial \theta'}$ must have a full rank (models not identified under the H_0 excluded)

Wald test

- Wald test is of the form

$$W = c(\hat{\theta})' (\hat{C}' \hat{V} \hat{C})^{-1} c(\hat{\theta}) \xrightarrow{d} \chi_Q^2$$

where g is the number of restrictions, \hat{V} is asymptotic variance of $\hat{\theta}$ and \hat{C} is the Jacobian matrix of $c(\hat{\theta})$ (matrix of first order derivatives)

- The idea of the test is that we reject H_0 if $c(\hat{\theta})$ differs significantly from zero
- Most important advantage of Wald test: only requires calculation of the estimate $\hat{\theta}$ for *unrestricted* model

- Practical limitations of this test:
 - $\hat{\theta}_0$ must be in the interior of Θ - not on the boundary
 - It is not invariant to how the nonlinear restrictions are imposed (in small samples)

Score test

- We impose the restriction directly on the model and estimate so called *restricted* model

$$\min_{\theta \in \Theta} \sum q(w_i, \theta) \text{ s.t. } c(\theta) = 0$$

- It is possible to prove that

$$LM = \left(\sum_{i=1}^N \tilde{s}_i \right)' \tilde{A}^{-1} \tilde{C}' \left(\tilde{C} \tilde{A}^{-1} \tilde{B}^{-1} \tilde{A} \tilde{C}' \right)^{-1} \tilde{C} \tilde{A}^{-1} \left(\sum_{i=1}^N \tilde{s}_i \right) / N \xrightarrow{d} \chi_Q^2$$

where all the matrices are calculated for restricted estimates of θ

- This is the heteroscedasticity robust form
- If *GIME* is true it is possible to simplify it to the form

$$LM = \left(\sum_{i=1}^N \tilde{\mathbf{s}}_i \right)' \tilde{\mathbf{M}}^{-1} \left(\sum_{i=1}^N \tilde{\mathbf{s}}_i \right)$$

- Where $\tilde{\mathbf{M}}$ can be estimated as $\sum_{i=1}^N \tilde{\mathbf{A}}_i$, $\sum_{i=1}^N \tilde{\mathbf{H}}_i$ or $\sum_{i=1}^N \tilde{\mathbf{s}}_i \tilde{\mathbf{s}}_i'$.
- The last estimate results in so called outer product of the score *LM* statistics which can be calculated as NR_o^2 from the regression of 1 on $\tilde{\mathbf{s}}_i$, where R_o^2 is uncentered R^2 .
- Wooldrige gives a similar way how to calculate the heteroscedasticity robust form of the test for *NLS*

- Vector $\sum_{i=1}^N \tilde{s}_i$ is called the score vector and is the first order derivative of the objective function calculated at *restricted* estimates
- The idea of the test is that if for the first order derivatives calculated for restricted estimates are far away from zero (f.o.c. for unrestricted model are far from being true), that we tend to reject H_0
- Advantages of *LM* test:
 - requires only the restricted estimates - ideal for specification testing
 - is invariant to specifications of the H_0
 - in some cases does not require that θ is in interior of Θ , leading case $\theta = (\theta_1, \theta_2)$ and $H_0 : \theta_2 = 0$

Tests based on the change in objective function

- The QLR test has the form

$$QLR = 2 \left[\sum_{i=1}^N q(\mathbf{w}_i, \tilde{\boldsymbol{\theta}}) - \sum_{i=1}^N q(\mathbf{w}_i, \hat{\boldsymbol{\theta}}) \right] \xrightarrow{d} \chi_Q^2$$

where $\tilde{\boldsymbol{\theta}}$ is the restricted estimate and $\hat{\boldsymbol{\theta}}$ is unrestricted estimate

- Intuition of this test is that if imposing restriction causes the big reduction in the objective function at the restricted maximum we tend to reject the H_0
- Advantages of this statistics:

- if the restricted and unrestricted model is easy to estimate then this statistic is very easy to calculate
- invariant to the specifications of the H_0
- Disadvantages
 - cannot be made robust to heteroscedasticity

Optimization methods

- For many problem is impossible to find analytically the maximum of the objective function $\sum_{i=1}^N q_i(\mathbf{w}_i, \boldsymbol{\theta})$
- In this case we find the maximum numerically
- One of the simplest methods the bisection method
- However it is not efficient if $\boldsymbol{\theta}$ has many elements
- In this cases we use so called gradient methods
- One of the most popular is so called Newton-Raphson method

Maximum likelihood estimation

- In most economic application we specify conditional probability given \mathbf{x}_i (we are not interested in distribution of exogenous variables)
- Then we use conditional maximum likelihood estimation *CMLE*
- We assume that the density is unknown up to the finite number of parameters (parametric estimation)

Example 3. (*Wooldridge*) Suppose that latent variable follows

$$y_i^* = \mathbf{x}_i \boldsymbol{\theta} + e_i$$

and $e_i \sim N(0, 1)$. We observe

$$y_i = \begin{cases} 1 & \text{if } y_i^* > 0 \\ 0 & \text{if } y_i^* \leq 0 \end{cases}$$

Probability

$$\begin{aligned} \Pr(y_i = 1 | \mathbf{x}_i) &= \Pr(\mathbf{x}_i \boldsymbol{\theta} + e_i > 0) = \Pr(e_i > -\mathbf{x}_i \boldsymbol{\theta}) \\ &= 1 - \Phi(-\mathbf{x}_i \boldsymbol{\theta}) = \Phi(\mathbf{x}_i \boldsymbol{\theta}) \end{aligned}$$

The probability of $\Pr(y_i = 0 | \mathbf{x}_i) = 1 - \Phi(\mathbf{x}_i \boldsymbol{\theta})$. The conditional density for an observation i can be written as

$$f(y_i | \mathbf{x}_i) = [\Phi(\mathbf{x}_i \boldsymbol{\theta})]^y [1 - \Phi(\mathbf{x}_i \boldsymbol{\theta})]^{1-y}$$

- The model is correctly specified model for conditional density if for some

$$\theta_0 \in \Theta$$

$$f(y|x;\theta) = p_0(y|x) \text{ for all } x \in \mathcal{X}$$

- Maximum likelihood estimation is special case of M-estimation as

$$E[\ell_i(\theta_0)|x_i] \geq E[\ell_i(\theta)|x_i] \text{ for all } \theta \in \Theta$$

where $\ell_i(\theta) = \ell_i(\mathbf{y}_i, \mathbf{x}_i, \theta) = \log[f(y_i|x_i;\theta)]$

- The M -estimator in the case of the Maximum Likelihood is found as the maximum of

$$\max_{\theta \in \Theta} \sum_{i=1}^N \ell_i(\theta)$$

Example 4. For probit model conditional log likelihood have the form $\ell_i(\theta) = y_i \log[\Phi(\mathbf{x}_i\theta)] + (1 - y_i) \log[1 - \Phi(\mathbf{x}_i\theta)]$

- Consistency and asymptotic normality of ML estimators are can be derived as a special case of the M -estimators
- The same is true about the test statistics (Wald, LM and LR)

Partial maximum likelihood methods

- Assume that we have the correct conditional density for each observation $f_t(y_t | \mathbf{x}_t, \boldsymbol{\theta}_o)$ - partial density
- However we do not assume that

$$p_o(\mathbf{y} | \mathbf{z}) = \prod_{i=1}^N f_t(y_t | \mathbf{x}_t, \boldsymbol{\theta}_o)$$

is the correct density function for all the sample given some set of variables \mathbf{z}

- It is often much simpler to define $f_t(y_t | \mathbf{x}_t, \boldsymbol{\theta}_o)$ than the functions which would fulfill the above requirement

- With the standard assumptions concerning M -estimators, it is possible to prove that the partial maximum likelihood estimator is consistent and asymptotically normal
- In order to obtain the correct variance matrix we have to use however the robust versions of variance estimators

Example 5. *(Probit with panel data). If the e_i are correlated then it is difficult to derive the conditional distribution of y_{it} given observation x_{it} and $x_{it-1}, \dots, x_{i1}, y_{it-1}, \dots, y_{i1}$ (y_{it} depends on them indirectly because of the correlation of e_i). However the conditional distribution of the probit model conditional only on x_{it} is the same as original one. Then the consistent estimator can be found using the pooled ML estimator and using the robust variance covariance matrix.*

- For panel data we can use the estimators designed for clusters. We define the same units as clusters and estimated covariance matrix is adjusted

properly.