

Máxima Verossimilhança

Bibliografia: Greene

Claudio Lucinda

FEA/USP



Overview

- 1 Máxima Verossimilhança
- 2 Propriedades do ML
- 3 A Igualdade da Matriz de Informação
- 4 Prova de Eficiência e Normalidade Assintótica
- 5 Estimando a Variância Assintótica do estimador ML



Máxima Verossimilhança

- O ponto de partida do Método da Máxima Verossimilhança é a imposição de uma pdf que assumimos que represente o processo gerador dos dados que estão na nossa amostra.
- Essa PDF é condicional a um vetor de parâmetros, e a densidade conjunta de n iid observações que seguem esta pdf é o produto das densidades individuais;

$$f(y_1, \dots, y_n | \theta) = \prod_{i=1}^n f(y_i | \theta) = L(\theta | \mathbf{y})$$

Essa pdf conjunta é chamada de Função Verossimilhança, e \mathbf{y} é a representação da amostra. Geralmente é mais simples trabalhar com o logaritmo da função verossimilhança, ou a log-verossimilhança:

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



Variáveis Independentes

- Como exemplo, vamos usar o modelo de Regressão Linear Clássico.
- Suponha que o termo erro seja normalmente distribuído
- Ou seja, condicional aos \mathbf{x}_i , y_i é normalmente distribuída com média $\mu_i = \mathbf{x}_i' \boldsymbol{\beta}$ e variância σ^2 .
- A verossimilhança fica sendo então:

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



Identificação em ML

Definition

O Vetor de Parâmetros θ é identificado (=estimável) se, para qualquer outro vetor $\theta^* \neq \theta$, para algum banco de dados \mathbf{y} , $L(\theta^* | \mathbf{y}) \neq L(\theta | \mathbf{y})$.



Propriedades do ML - Condições de Regularidade

- R1. As primeiras três derivadas de $\ln f(y_i | \theta)$ com respeito a θ são finitas e contínuas para quase todos os y_i e para todos os θ . Isso assegura a existência de uma aproximação de série de Taylor e garante uma variância finita das derivadas de $\ln L$.
- R2. As condições necessárias para se obter as esperanças das primeiras e segundas derivadas de $\ln f(y_i | \theta)$ são atendidas.
- R3. Para todos os valores de θ , $|\partial^3 \ln f(y_i | \theta) / \partial \theta_j \partial \theta_k \partial \theta_l|$ é menor que uma função com esperança finita. Isso vai permitir que a série de Taylor seja truncada no segundo termo.



Propriedades do ML

- M1. Consistência: $\text{plim } \hat{\theta} = \theta_0$.
- M2. Normalidade Assintótica: $\hat{\theta} \stackrel{a}{\sim} N \left[\theta_0, \{I(\theta_0)\}^{-1} \right]$, em que

$$I(\theta_0) = -E_0 \left[\partial^2 \ln L / \partial \theta_0 \partial \theta_0' \right].$$

- M3. Eficiência Assintótica: $\hat{\theta}$ é assintoticamente eficiente e alcança o limite inferior de Cramér-Rao para estimadores consistentes.
- M4. Invariância: O estimador ML de $\gamma_0 = \mathbf{c}(\theta_0)$ is $\mathbf{c}(\hat{\theta})$ if $\mathbf{c}(\theta_0)$ é uma função contínua e continuamente diferenciável.



Momentos das Derivadas do ML

- D1. $\ln f(y_i | \theta)$, $\mathbf{g}_i = \partial \ln f(y_i | \theta) / \partial \theta$, and $\mathbf{H}_i = \partial^2 \ln f(y_i | \theta) / \partial \theta \partial \theta'$, $i = 1, \dots, n$, todas são amostras aleatórias de variáveis aleatórias (porque assumimos amostragem aleatória). A notação $\mathbf{g}_i(\theta_0)$ and $\mathbf{H}_i(\theta_0)$ indica a derivada avaliada em θ_0 .
- D2. $E_0[\mathbf{g}_i(\theta_0)] = \mathbf{0}$.
- D3. $\text{Var}[\mathbf{g}_i(\theta_0)] = -E[\mathbf{H}_i(\theta_0)]$.

Equação da Verossimilhança

A função verossimilhança é

$$\ln L(\boldsymbol{\theta} \mid \mathbf{y}) = \sum_{i=1}^n \ln f(y_i \mid \boldsymbol{\theta})$$

O vetor de primeiras derivadas, ou vetor de score, é dado por

$$\mathbf{g} = \frac{\partial \ln L(\boldsymbol{\theta} \mid \mathbf{y})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n \frac{\partial \ln f(y_i \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=1}^n \mathbf{g}_i$$

Como estamos apenas somando termos, decorre das propriedades D1 e D2 do slide anterior que em $\boldsymbol{\theta}_0$,

$$E_0 \left[\frac{\partial \ln L(\boldsymbol{\theta}_0 \mid \mathbf{y})}{\partial \boldsymbol{\theta}_0} \right] = E_0 [\mathbf{g}_0] = \mathbf{0}.$$



A Igualdade da Matriz de Informação - I

O Hessiano da log verossimilhança é

$$\mathbf{H} = \frac{\partial^2 \ln L(\boldsymbol{\theta} \mid \mathbf{y})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} = \sum_{i=1}^n \frac{\partial^2 \ln f(y_i \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} = \sum_{i=1}^N \mathbf{H}_i.$$

Avaliando em $\boldsymbol{\theta}_0$ e definindo

$$E_0 [\mathbf{g}_0 \mathbf{g}_0'] = E_0 \left[\sum_{i=1}^n \sum_{j=1}^n \mathbf{g}_{0i} \mathbf{g}_{0j}' \right]$$

A Igualdade da Matriz de Informação – II

Devido a D1, podemos tirar os termos com subscritos diferentes e chegar em

$$E_0 [\mathbf{g}_0 \mathbf{g}_0'] = E_0 \left[\sum_{i=1}^n \mathbf{g}_{0i} \mathbf{g}_{0i}' \right] = E_0 \left[\sum_{i=1}^n (-\mathbf{H}_{0i}) \right] = -E_0 [\mathbf{H}_0]$$

Tal que

$$\begin{aligned} \text{Var}_0 \left[\frac{\partial \ln L(\boldsymbol{\theta}_0 | \mathbf{y})}{\partial \boldsymbol{\theta}_0} \right] &= E_0 \left[\left(\frac{\partial \ln L(\boldsymbol{\theta}_0 | \mathbf{y})}{\partial \boldsymbol{\theta}_0} \right) \left(\frac{\partial \ln L(\boldsymbol{\theta}_0 | \mathbf{y})}{\partial \boldsymbol{\theta}_0'} \right) \right] \\ &= -E_0 \left[\frac{\partial^2 \ln L(\boldsymbol{\theta}_0 | \mathbf{y})}{\partial \boldsymbol{\theta}_0 \partial \boldsymbol{\theta}_0'} \right]. \end{aligned}$$

Normalidade Assintótica

No estimador ML, o gradiente da função verossimilhança é igual a zero por construção, então

$$\mathbf{g}(\hat{\theta}) = \mathbf{0}$$

Fazendo uma expansão de Taylor de segunda ordem em volta dos valores verdadeiros θ_0 . Vamos usar o teorema do valor médio para truncar a expansão no segundo termo.

$$\mathbf{g}(\hat{\theta}) = \mathbf{g}(\theta_0) + \mathbf{H}(\bar{\theta}) (\hat{\theta} - \theta_0) = \mathbf{0}.$$



Normalidade Assintótica II

O Hessiano é avaliado em um ponto $\bar{\theta}$ entre $\hat{\theta}$ e θ_0 ($\bar{\theta} = w\hat{\theta} + (1-w)\theta_0$ para algum $0 < w < 1$). Reorganizando a função e multiplicando o resultado por \sqrt{n} :

$$\sqrt{n}(\hat{\theta} - \theta_0) = [-\mathbf{H}(\bar{\theta})]^{-1} [\sqrt{n}\mathbf{g}(\theta_0)]$$

Uma vez que $\text{plim}(\hat{\theta} - \theta_0) = \mathbf{0}$, $\text{plim}(\hat{\theta} - \bar{\theta}) = 0$ também. As derivadas segundas são funções contínuas.



Normalidade Assintótica III

Se existir uma distribuição limite, ela seria

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} [-\mathbf{H}(\boldsymbol{\theta}_0)]^{-1} [\sqrt{n}\mathbf{g}(\boldsymbol{\theta}_0)]$$

Dividindo $\mathbf{H}(\boldsymbol{\theta}_0)$ e $\mathbf{g}(\boldsymbol{\theta}_0)$ por n , temos

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} \left[-\frac{1}{n}\mathbf{H}(\boldsymbol{\theta}_0)\right]^{-1} [\sqrt{n}\bar{\mathbf{g}}(\boldsymbol{\theta}_0)]$$



Normalidade Assintótica IV

Aplicando o Teorema Central do Limite de Lindberg-Lévy a $[\sqrt{n}\bar{\mathbf{g}}(\theta_0)]$, uma vez que é \sqrt{n} vezes a média de uma amostra aleatória.

A variância limite de $[\sqrt{n}\bar{\mathbf{g}}(\theta_0)]$ is $-E_0[(1/n)\mathbf{H}(\theta_0)]$, so

$$\sqrt{n}\bar{\mathbf{g}}(\theta_0) \xrightarrow{d} N\left\{\mathbf{0}, -E_0\left[\frac{1}{n}\mathbf{H}(\theta_0)\right]\right\}$$

Pelo Teorema de Chebyshev (Teorema D.2 do Greene),
 $\text{plim}[-(1/n)\mathbf{H}(\theta_0)] = -E_0[(1/n)\mathbf{H}(\theta_0)]$. Como é uma matriz constante, podemos reorganizar

$$\left[-\frac{1}{n}\mathbf{H}(\theta_0)\right]^{-1} \sqrt{n}\bar{\mathbf{g}}(\theta_0) \xrightarrow{d} N\left[\mathbf{0}, \left\{-E_0\left[\frac{1}{n}\mathbf{H}(\theta_0)\right]\right\}^{-1} \left\{-E_0\left[\frac{1}{n}\mathbf{H}(\theta_0)\right]\right\} \left\{-E_0\left[\frac{1}{n}\mathbf{H}(\theta_0)\right]\right\}^{-1}\right]$$



Normalidade Assintótica V

Voltando

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N\left[\mathbf{0}, \left\{-E_0\left[\frac{1}{n}\mathbf{H}(\boldsymbol{\theta}_0)\right]\right\}^{-1}\right]$$

Que nos dá a distribuição assintótica do MLE:

$$\hat{\boldsymbol{\theta}} \overset{a}{\sim} N\left[\boldsymbol{\theta}_0, \{\mathbf{I}(\boldsymbol{\theta}_0)\}^{-1}\right]$$



Eficiência Assintótica e o Cramér-Rao Lower Bound

Theorem

Supondo que a densidade de y_i atenda às condições de regularidade R1 – R3, a variância assintótica de um estimador normalmente distribuído consistente e assintoticamente normalmente distribuído θ_0 será sempre pelo menos tão grande quanto:

$$[\mathbf{I}(\theta_0)]^{-1} = \left(-E_0 \left[\frac{\partial^2 \ln L(\theta_0)}{\partial \theta_0 \partial \theta_0'} \right] \right)^{-1} = \left(E_0 \left[\left(\frac{\partial \ln L(\theta_0)}{\partial \theta_0} \right) \left(\frac{\partial \ln L(\theta_0)}{\partial \theta_0} \right)' \right] \right)^{-1}.$$



Variância Assintótica do ML

Se a forma dos valores esperados das derivadas segundas da log verossimilhança é conhecida, então podemos avaliar esta fórmula em $\hat{\theta}$ e encontrar a matriz VC do ML:

$$[\mathbf{I}(\theta_0)]^{-1} = \left\{ -E_0 \left[\frac{\partial^2 \ln L(\theta_0)}{\partial \theta_0 \partial \theta_0'} \right] \right\}^{-1}$$

Como isso quase nunca é fácil, existem alternativas. A primeira delas é

$$[\hat{\mathbf{I}}(\hat{\theta})]^{-1} = \left(-\frac{\partial^2 \ln L(\hat{\theta})}{\partial \hat{\theta} \partial \hat{\theta}'} \right)^{-1}.$$

Ou seja, avaliando as derivadas segundas na função verossimilhança em torno das estimativas de ML.



Variância Assintótica do Estimador ML

Outro estimador é baseado no resultado abaixo:

$$[\hat{\mathbf{l}}(\hat{\boldsymbol{\theta}})]^{-1} = \left[\sum_{i=1}^n \hat{\mathbf{g}}_i \hat{\mathbf{g}}_i' \right]^{-1} = [\hat{\mathbf{G}}' \hat{\mathbf{G}}]^{-1},$$

em que

$$\hat{\mathbf{g}}_i = \frac{\partial \ln f(\mathbf{x}_i, \hat{\boldsymbol{\theta}})}{\partial \hat{\boldsymbol{\theta}}}$$

e

$$\hat{\mathbf{G}} = [\hat{\mathbf{g}}_1, \hat{\mathbf{g}}_2, \dots, \hat{\mathbf{g}}_n]'$$

