

Lecture 1: Maximum Likelihood Estimator



Professor: Mauricio Sarrias

Universidad de Talca

2020

- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix

- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality

- 3 Estimation of Variance

- 4 Testing
 - Intuition
 - The Trinity
 - Proof, Proof and More Proof

Reading

Reading (Mandatory):

- (Ruud)- Chapters 14 and 15.
- Buse, A. (1982). The likelihood ratio, Wald, and Lagrange multiplier tests: An expository note. *The American Statistician*, 36(3a), 153-157.

Suggested:

- (Winkelmann & Boes)- Chapters 2 and 3

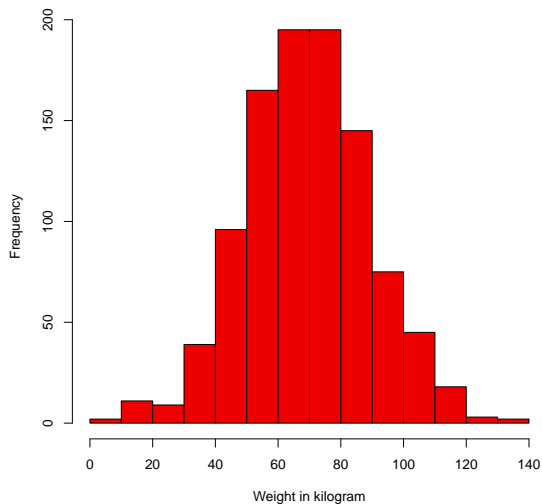
Goals

- Understand the logic behind the Maximum Likelihood Estimator.
- Derive the asymptotic properties of the MLE.
- Understand and derive the basic test for the MLE.

- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix
- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality
- 3 Estimation of Variance
- 4 Testing
 - Intuition
 - The Trinity
 - Proof, Proof and More Proof

Motivating Example I

Let's assume we weighed 1000 people from Talca



Motivating Example I

The goal of maximum likelihood is to find the optimal way to fit a distribution to the data.

Remark I

Generally, we can write the probability or density function of $y_i = 1, \dots, n$ as $f(y_i; \theta)$, where y_i is the i th draw from the population and θ is the parameter of the distribution.

Remark II

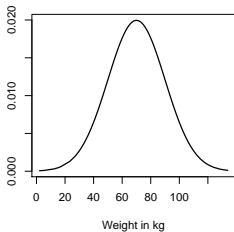
We usually assume independent sampling, i.e., the i th draw from the population is independent from all other draws $i' \neq i$

So the question is:

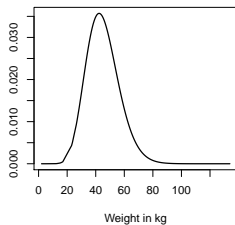
Which distribution does fit the previous weight data?

Motivating Example I

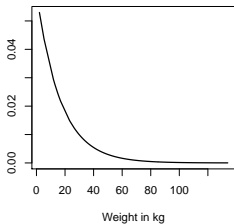
Normal Distribution



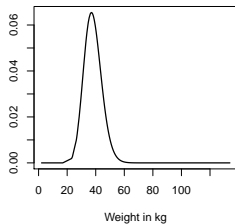
Chi-squared Distribution



Exponential Distribution



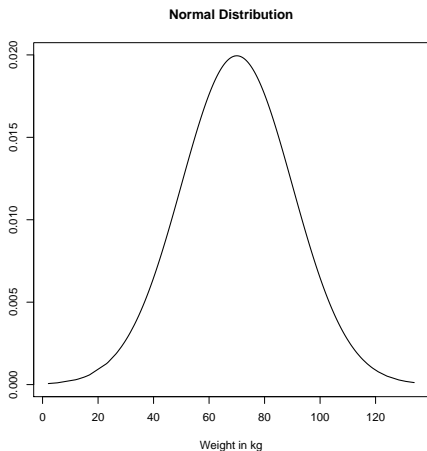
Gamma Distribution



Motivating Example I

It seems that the normal distribution is the best option.

- We expect most of the weights to be close to the mean.
- We expect the weights to be relatively symmetrical around the mean.
- Ok ..., but not every normal fits the our data.
- What mean, μ , and variance, σ^2 , are the best “estimates”?

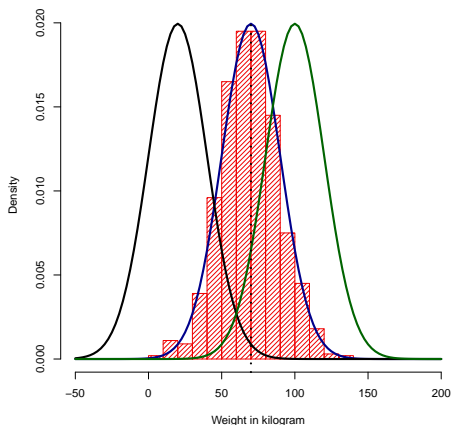


Motivating Example I

Maximum Likelihood Principle

- 1 We observe some data.
- 2 We pick the distribution we think generated the data.
- 3 We find the estimator(s) of the distribution, $\hat{\theta}$, that makes more likely the sample we are observing.

IOW, the problem consists on estimating an unknown parameter of a population when the population distribution is known (up to the unknown parameter)



Motivating Example II

Example

A random sample of 100 trials was performed and 10 resulted in success. What can be inferred about the unknown probability of success p_0 ?

Note that we are observing the sample; somehow we know the distribution; and we are asking what is \hat{p} that makes more likely the sample we are observing.

Motivating Example II

For any potential value of p for the probability of success, the probability of y successes from n trials is given by:

$$f(y; n, p) = \Pr(Y = y) = \binom{n}{y} p^y (1 - p)^{n-y}$$

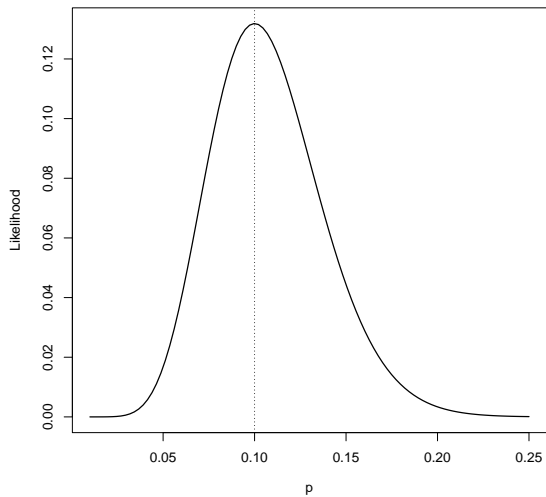
where

$$\binom{n}{y} = \frac{n!}{y!(n-y)!}$$

With $y = 10$ successes from $n = 100$ trials,

$$\begin{aligned} L(p) &= \Pr(Y = 10) \\ &= \frac{100!}{90!10!} p^{10} (1 - p)^{90} \\ &= 1.731 \times 10^{13} \times p^{10} (1 - p)^{90} \end{aligned}$$

Motivating Example II



Likelihood Function

The **likelihood function** denoted by capital L is:

$$L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}) = \prod_{i=1}^n L(\boldsymbol{\theta}; y_i|\mathbf{x}_i) = \prod_{i=1}^n f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$$

where $\mathbf{y} = (y_1, \dots, y_n)$.

- $L(\boldsymbol{\theta}; y_i|\mathbf{x}_i)$ is the likelihood contribution of the **i -th observation**,
 - $L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$ is the likelihood function of the **whole sample**.
- The likelihood function says that, for any given sample $\mathbf{y}|\mathbf{X}$, the likelihood of having obtained that particular sample depends on the parameter $\boldsymbol{\theta}$.
- Whenever we can write down the **joint probability function** of the sample we can in principle use ML estimation.

Log Likelihood Function

The **log-likelihood function** is:

$$\ln L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X}) = \ln L(\boldsymbol{\theta}) = \ln \underbrace{\left(\prod_{i=1}^n f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \right)}_{f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})} = \sum_{i=1}^n \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$$

- The log-likelihood function is a **monotonically increasing** function of $L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$:
 - ▶ Any maximizing value $\hat{\boldsymbol{\theta}}$ of $\ln L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$ must also maximize $L(\boldsymbol{\theta}, \mathbf{y}|\mathbf{X})$.
- Taking logarithms which converts products into sums.
 - ▶ It allows some simplification in the numerical determination of the MLE.
 - ▶ Likelihood values are often extremely small (but can also be extremely large). Numerical optimization of the likelihood highly problematic.
 - ▶ Simplification of the study of the properties of the estimator.

Example

Example (Binomial Example)

Let $\{Y_n\}$ be a random sample of a binomial r.v with parameters (n, p) , where n is assumed to be known and p unknown. The likelihood function for individual i is given by:

$$L_i(p; y_i) = f(y_i; p) = \binom{n}{y_i} p^{y_i} (1 - p)^{n - y_i}$$

Since the sample is iid, the likelihood function:

$$L(p; \mathbf{y}) = f(y_1, y_2, \dots, y_n; p) = \prod_{i=1}^n f(y_i; p) = \prod_{i=1}^n \binom{n}{y_i} p^{y_i} (1 - p)^{n - y_i}$$

Example

Example (Binomial Example)

Taking the log we get the log-likelihood function:

$$\begin{aligned}\ln L(p; \mathbf{y}) &= \ln \left(\prod_{i=1}^n f(y_i; p) \right) \\ &= \ln \left(\prod_{i=1}^n \binom{n}{y_i} p^{y_i} (1-p)^{(n-y_i)} \right) \\ &= \ln \left[p^{\sum_{i=1}^n y_i} (1-p)^{(n^2 - \sum_{i=1}^n y_i)} \prod_{i=1}^n \binom{n}{y_i} \right] \\ &= \ln \prod_{i=1}^n \binom{n}{y_i} + \left(\sum_{i=1}^n y_i \right) \ln p + \left(n^2 - \sum_{i=1}^n y_i \right) \ln(1-p)\end{aligned}$$

Example: Linear Regression

Example (Linear Regression)

Consider that $\{y_i, \mathbf{x}_i\}$ is i.i.d, and $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i$, where $\epsilon_i | \mathbf{x}_i \sim N(0, \sigma_0^2)$. So, with $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$ and $\mathbf{w}_i = (y_i, \mathbf{x}_i^\top)^\top$, the conditional pdf is

$$\begin{aligned} f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) &= \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp \left[-\frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_0)^2}{2\sigma_0^2} \right] \\ &= \phi(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_0, \sigma_0^2) \end{aligned}$$

The joint pdf of the sample is:

$$\begin{aligned} \prod_{i=1}^n f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) &= [2\pi\sigma_0^2]^{n/2} \exp \left[-\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)}{2\sigma_0^2} \right] \\ &= \phi(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0, \sigma_0^2 \cdot \mathbf{I}_n) \end{aligned}$$

The parameter space is Θ is $\mathbb{R}^K \times \mathbb{R}_{++}$, where K is the dimension of $\boldsymbol{\beta}$ and \mathbb{R}_{++} is the set of positive real numbers reflecting the a priori restriction that $\sigma_0^2 > 0$

1 Introduction

- Motivation
- **Maximum Likelihood Estimator**
- Identification
- The Score Function
- The Information Matrix

2 Asymptotic Properties

- Consistency
- Asymptotic Normality

3 Estimation of Variance

4 Testing

- Intuition
- The Trinity
- Proof, Proof and More Proof

Maximum Likelihood Estimator

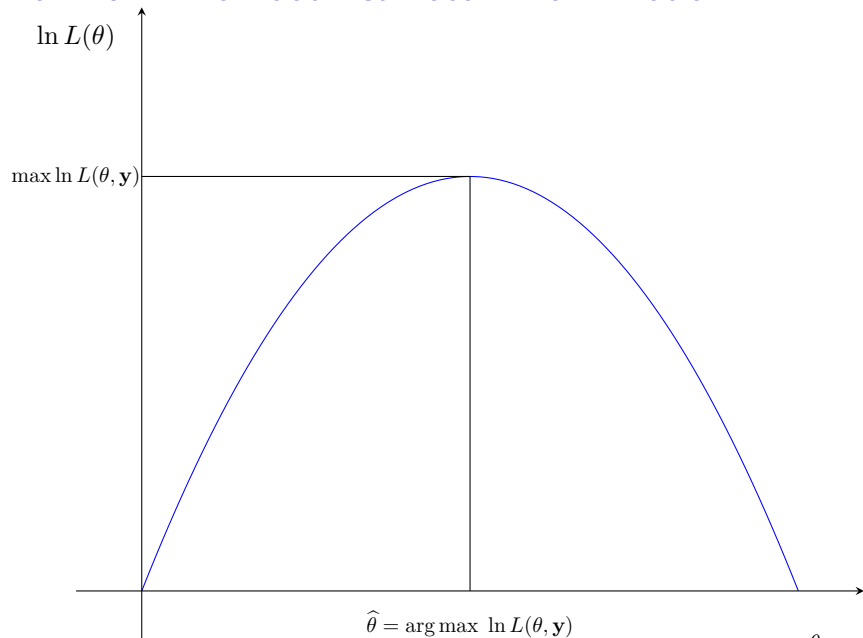
Definition (ML Estimator)

The MLE is a value of the parameter vector that maximizes the sample average log-likelihood function:

$$\hat{\boldsymbol{\theta}}_n \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$$

where Θ denotes the parameter space in which the parameter vector $\boldsymbol{\theta}$ lies. Usually $\Theta = \mathbb{R}^K$.

Maximum Likelihood Estimator: Maximization



Maximum Likelihood Estimator

Remark:

By the nature of the objective function, the MLE is the **estimator which makes the observed data most likely to occur**. In other words, the MLE is the best “rationalization” of what we observed.

Population analogous

$$\mathbb{E} [\ln L(\boldsymbol{\theta}; \mathbf{y}|\mathbf{X})] \equiv \int \ln L(\boldsymbol{\theta}; \mathbf{y}|\mathbf{X}) dF(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}_0)$$

where $F(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}_0)$ is the joint CDF of (\mathbf{y}, \mathbf{X})

Maximum Likelihood Estimator

Assumption I: Distribution

The sample $\{y_i, \mathbf{x}_i\}$ is i.i.d with **true** conditional density $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$.

Since the sample is i.i.d, we can write:

$$f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) = f(y_1|\mathbf{x}_1; \boldsymbol{\theta}) \times f(y_2|\mathbf{x}_2; \boldsymbol{\theta}) \times \dots \times f(y_n|\mathbf{x}_n; \boldsymbol{\theta})$$

Expected Log-Likelihood Inequality

Is $\mathbb{E} [\ln L(\boldsymbol{\theta}; \mathbf{y}|\mathbf{X})]$ maximized at $\boldsymbol{\theta}_0$?

Assumption II: Dominance I

$\mathbb{E} [\sup_{\boldsymbol{\theta} \in \Theta} |\ln L(\boldsymbol{\theta}; \mathbf{y}|\mathbf{X})|]$ exists.

Lemma (Expected Log-likelihood Inequality)

If Dominance I assumption holds, then

$$\mathbb{E} [\ln f(y|\mathbf{x}; \boldsymbol{\theta})] \leq \mathbb{E} [\ln f(y|\mathbf{x}; \boldsymbol{\theta}_0)]$$

Example

The conditional log-likelihood function of $y_i | \mathbf{x}_i \sim N(\mathbf{x}_i^\top \boldsymbol{\beta}_0, \sigma^2)$ is

$$\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = -0.5 \log(2\pi\sigma^2) - \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{2\sigma^2} \quad (1)$$

The conditional expectation is:

$$\mathbb{E}[\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i] = -0.5 \log(2\pi\sigma^2) - \frac{\mathbb{E}[(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2]}{2\sigma^2}$$

Note that:

$$\begin{aligned} \mathbb{E}[(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2] &= \mathbb{E}[(\mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2] \quad \because y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + \epsilon_i \\ &= \mathbb{E}[(\epsilon_i + \mathbf{x}_i^\top (\boldsymbol{\beta}_0 - \boldsymbol{\beta}))^2] \\ &= \mathbb{E}[\epsilon_i^2 + 2\epsilon_i \mathbf{x}_i^\top (\boldsymbol{\beta}_0 - \boldsymbol{\beta}) + (\boldsymbol{\beta}_0 - \boldsymbol{\beta})^\top \mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}_0 - \boldsymbol{\beta})] \\ &= \mathbb{E}(\epsilon_i^2) + (\boldsymbol{\beta}_0 - \boldsymbol{\beta})^\top \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top) (\boldsymbol{\beta}_0 - \boldsymbol{\beta}) \\ &= \sigma_0^2 + (\boldsymbol{\beta}_0 - \boldsymbol{\beta})^\top \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top) (\boldsymbol{\beta}_0 - \boldsymbol{\beta}) \end{aligned}$$

When is this expectation finite?

Example

The last term is finite if $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top)$ is. This implies that \mathbf{X} is full-column rank.

$$\mathbb{E}[\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i] = -0.5 \log(2\pi\sigma^2) - \frac{\sigma_0^2 + (\boldsymbol{\beta}_0 - \boldsymbol{\beta})^\top \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top) (\boldsymbol{\beta}_0 - \boldsymbol{\beta})}{2\sigma^2}$$

Now, $\mathbb{E}[\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i]$ is uniquely maximized at $\mathbf{x}_i^\top \boldsymbol{\beta} = \mathbf{x}_i^\top \boldsymbol{\beta}_0$ and $\sigma^2 = \sigma_0^2$

1 Introduction

- Motivation
- Maximum Likelihood Estimator
- **Identification**
- The Score Function
- The Information Matrix

2 Asymptotic Properties

- Consistency
- Asymptotic Normality

3 Estimation of Variance

4 Testing

- Intuition
- The Trinity
- Proof, Proof and More Proof

Identification

- Before employing MLE, it is necessary to check whether the data-generating process is sufficiently informative about the parameters of the model.
- Recall OLS: $\hat{\beta}$ to be unique \mathbf{X} must be full-column rank. Otherwise, ...
- The question is: is the population $\mathbb{E}[\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$ uniquely maximized at $\boldsymbol{\theta}_0$?
 - ▶ If there exists another $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ that maximized $\mathbb{E}[\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$, then MLE is not identified.
- This is satisfied if (**conditional density identification**):

$$f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \neq f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0) \quad \forall \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$$

Identification

Definition (Global Identification)

The parameter vector θ_0 is globally identified in Θ if, for every $\theta_1 \in \Theta$, $\theta \neq \theta_1$ implies that:

$$\Pr [f(y_i|\mathbf{x}_i; \theta_0) \neq f(y_i|\mathbf{x}_i; \theta_1)] > 0$$

Assumption III: Global Identification

Every parameter vector $\theta_0 \in \Theta$ is globally identified.

Identification

Lemma (Strict Expected Log-Likelihood Inequality)

*Under the Assumptions of **Distribution**, **Dominance I** and **Global Identification**, then*

$$\boldsymbol{\theta} \neq \boldsymbol{\theta}_0 \implies \mathbb{E} [\ln f(y|\mathbf{x}; \boldsymbol{\theta})] < \mathbb{E} [\ln f(y|\mathbf{x}; \boldsymbol{\theta}_0)]$$

In words, the expected value of the log-likelihood is maximized at the true value of the parameters.

Proof.

Let $\mathbf{w} = (y, \mathbf{x}')'$ and define

$$a(\mathbf{w}) \equiv f(y|\mathbf{x}; \boldsymbol{\theta})/f(y|\mathbf{x}; \boldsymbol{\theta}_0)$$

First, WTS that $a(\mathbf{w}) \neq 1$ with positive probability, so that $a(\mathbf{w})$ is nonconstant random variable (so, we can apply Jensen's Inequality).

$$\begin{aligned} a(\mathbf{w}) \neq 1 &\iff f(y|\mathbf{x}; \boldsymbol{\theta}) \neq f(y|\mathbf{x}; \boldsymbol{\theta}_0) \\ \Pr[a(\mathbf{w}) \neq 1] &\iff \Pr[f(y|\mathbf{x}; \boldsymbol{\theta}) \neq f(y|\mathbf{x}; \boldsymbol{\theta}_0)] \end{aligned}$$

But, by Global Identification:

$$\Pr[f(y|\mathbf{x}; \boldsymbol{\theta}) \neq f(y|\mathbf{x}; \boldsymbol{\theta}_0)] > 0 \implies \Pr[a(\mathbf{w}) \neq 1] > 0$$

Now, WTS $\mathbb{E}[\log a(\mathbf{w})] < \log \{\mathbb{E}[a(\mathbf{w})]\}$. We use the strict version of **Jensen's inequality** which states that if $c(x)$ is a strictly concave function and x is nonconstant random variable, then $\mathbb{E}[c(x)] < c[\mathbb{E}(x)]$ □

Proof.

Set $c(x) = \log(x)$, since $\log(x)$ is strictly concave and $a(\mathbf{w})$ is non-constant. Therefore

$$\mathbb{E}[\log a(\mathbf{w})] < \log \{\mathbb{E}[a(\mathbf{w})]\}$$

Now, WTS that $\mathbb{E}(a(\mathbf{w})) = 1$. Note that the conditional mean of $a(\mathbf{w})$ equals 1 because:

$$\begin{aligned}\mathbb{E}[a(\mathbf{w})|\mathbf{x}] &= \int a(\mathbf{w})f(y|\mathbf{x};\boldsymbol{\theta}_0)dy \\ &= \int \frac{f(y|\mathbf{x};\boldsymbol{\theta})}{f(y|\mathbf{x};\boldsymbol{\theta}_0)}f(y|\mathbf{x};\boldsymbol{\theta}_0)dy \\ &= \int f(y|\mathbf{x};\boldsymbol{\theta})dy \\ &= 1\end{aligned}$$

By the Law of Total Expectations $\mathbb{E}[a(\mathbf{w})] = 1$. Combining the results:

$$\mathbb{E}[\log(a(\mathbf{w}))] < \log(1) = 0$$

But $\log(a(\mathbf{w})) = \log f(y|\mathbf{x};\boldsymbol{\theta}) - \log f(y|\mathbf{x};\boldsymbol{\theta}_0)$. □

1 Introduction

- Motivation
- Maximum Likelihood Estimator
- Identification
- **The Score Function**
- The Information Matrix

2 Asymptotic Properties

- Consistency
- Asymptotic Normality

3 Estimation of Variance

4 Testing

- Intuition
- The Trinity
- Proof, Proof and More Proof

Differentiability

Assumption IV: Integrability

The pdf $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ is twice continuously differentiable in $\boldsymbol{\theta}$ for all $\boldsymbol{\theta} \in \Theta$. Furthermore, the support $\mathcal{S}(\boldsymbol{\theta})$ of $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ does not depend on $\boldsymbol{\theta}$, and differentiation and integration are interchangeable in the sense that

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \int_{\mathcal{S}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) &= \int_{\mathcal{S}} \frac{\partial}{\partial \boldsymbol{\theta}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \\ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \int_{\mathcal{S}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) &= \int_{\mathcal{S}} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta})\end{aligned}$$

and

$$\begin{aligned}\frac{\partial \mathbb{E} [\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i = x_i]}{\partial \boldsymbol{\theta}} &= \mathbb{E} \left[\left. \frac{\partial \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right| \mathbf{x}_i = x_i \right] \\ \frac{\partial^2 \mathbb{E} [\ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) | \mathbf{x}_i = x_i]}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} &= \mathbb{E} \left[\left. \frac{\partial^2 \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right| \mathbf{x}_i = x_i \right]\end{aligned}$$

where all terms exists. In this case, we denote the support of $F(y)$ simply by \mathcal{S} .

The Score Function

Definition (Score Function)

The score function is defined as the vector of first partial derivatives of the log-likelihood function with respect to the parameter vector $\boldsymbol{\theta}$:

$$\mathbf{s}(\mathbf{w}, \boldsymbol{\theta}) = \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \begin{pmatrix} \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1} \\ \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2} \\ \vdots \\ \frac{\partial \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_K} \end{pmatrix}$$

The score vector for observation i is:

$$\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) = \frac{\partial \ln f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

Because of the additivity of terms in the log-likelihood function, we can write:

$$\mathbf{s}(\mathbf{w}, \boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})$$

Score Identity

Lemma (Score Identity)

Under *Integrability and Distribution Assumption*:

$$\mathbb{E}[\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})] = \mathbf{0}$$

- We have to be clear whether we are speaking about the score of a single observation $\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})$ or the score of the sample $\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})$.
- Since under random sampling, $\mathbf{s}(\mathbf{w}, \boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})$, it is sufficient to establish that $\mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})] = \mathbf{0}$

Proof.

First, we derive an integral property of pdf. Because we are **assuming** $F(y|x; \boldsymbol{\theta})$ is a proper cdf.,

$$\int_{\mathcal{S}} dF(y_i|\mathbf{x}_i; \boldsymbol{\theta}) = \int_{\mathcal{S}} f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i = 1 \quad (2)$$

$\forall \boldsymbol{\theta} \in \boldsymbol{\Theta}$. Given **differentiability**, we can differentiate both sides of this equality with respect to $\boldsymbol{\theta}$:

$$\mathbf{0} = \int_{\mathcal{S}} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i \quad (3)$$

This equation states how changes in $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ resulting from changes in $\boldsymbol{\theta}$ are restricted by (2). We can rewrite (3) as

$$\begin{aligned} \mathbf{0} &= \int_{\mathcal{S}} \frac{f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{f(y_i|\mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i \\ \mathbf{0} &= \int_{\mathcal{S}} \frac{1}{f(y_i|\mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial f(y_i|\mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \underbrace{dF(y_i|\mathbf{x}_i; \boldsymbol{\theta})}_{f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) dy_i} \end{aligned} \quad (4)$$

Proof.

Now we interpret this integral equation as an expectation. Consider:

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\theta}} \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &\equiv \frac{1}{f(y_i | \mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \\ \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) &\equiv \frac{1}{f(y_i | \mathbf{x}_i; \boldsymbol{\theta})} \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \\ \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &\equiv \frac{\partial}{\partial \boldsymbol{\theta}} f(y_i | \mathbf{x}_i; \boldsymbol{\theta})\end{aligned}\tag{5}$$

Then, substituting into (4)

$$\int_{\mathcal{S}} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) dF(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \mathbf{0}$$

This holds for any $\boldsymbol{\theta} \in \boldsymbol{\Theta}$, in particular, for $\boldsymbol{\theta} = \boldsymbol{\theta}_0$. Setting $\boldsymbol{\theta} = \boldsymbol{\theta}_0$, we obtain:

$$\int_{\mathcal{S}} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) dF(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) = \mathbf{0}$$

$$\int_{\mathcal{S}} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) dF(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0) = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) | \mathbf{x}] = \mathbf{0}$$

Then, by Law of Total Expectations, we obtain the desired result. □

What if the support depend on θ ?

In this case the support is $\mathcal{S}(\theta) = A(\theta) \leq y \leq B(\theta)$. By definition:

$$\int_{A(\theta)}^{B(\theta)} f(y|x; \theta) dy = 1$$

Now, using the Leibnitz's theorem gives:

$$\frac{\partial \int_{A(\theta)}^{B(\theta)} f(y|x; \theta) dy}{\partial \theta} = 0$$
$$\int_{A(\theta)}^{B(\theta)} \frac{\partial f(y|x; \theta)}{\partial \theta} dy + f(B(\theta)|\theta) \frac{\partial B(\theta)}{\partial \theta} - f(A(\theta)|\theta) \frac{\partial A(\theta)}{\partial \theta} = 0$$

To interchange the operations of differentiation and integration we need the second and third terms go to zero. The **necessary condition** is that

$$\lim_{y \rightarrow A(\theta)} f(y|x; \theta) = 0$$

$$\lim_{y \rightarrow B(\theta)} f(y|x; \theta) = 0$$

Sufficient conditions are that the support does not depend on the parameter, which means that $\partial A(\theta)/\partial \theta = \partial B(\theta)/\partial \theta = 0$ or that the density is zero at the terminal points.

1 Introduction

- Motivation
- Maximum Likelihood Estimator
- Identification
- The Score Function
- **The Information Matrix**

2 Asymptotic Properties

- Consistency
- Asymptotic Normality

3 Estimation of Variance

4 Testing

- Intuition
- The Trinity
- Proof, Proof and More Proof

Hessian

- Since we are doing an optimization analysis, we need the Hessian Matrix.

$$\mathbf{H}(\mathbf{w}; \boldsymbol{\theta}) = \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\tau} = \begin{pmatrix} \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1^2} & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_1 \partial \theta_K} \\ \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_2 \partial \theta_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_K \partial \theta_1} & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_K \partial \theta_K} & \cdots & \frac{\partial^2 \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta})}{\partial \theta_K^2} \end{pmatrix}$$

If the log-likelihood function is concave in $\boldsymbol{\theta}$, $\mathbf{H}(\mathbf{w}; \boldsymbol{\theta})$ is said to be negative definite. In the scalar case, for $K = 1$, this simply means that the second derivative of the log-likelihood function is negative.

Hessian

Because of the additivity of terms in the log-likelihood function:

$$\mathbf{H}(\mathbf{w}; \boldsymbol{\theta}) = \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) \quad \text{where} \quad \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) = \frac{\partial^2 \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}$$

Remark

It is important to keep in mind that both the score and Hessian depend on the sample and are therefore random variables (they differ in repeated samples).

Information Identity

- To analyze the variance and the limiting distribution of the ML estimator, we require some results on the **Fisher information matrix**.
- It is very related to the Hessian matrix.
- The information matrix of a sample is simply defined as the negative expectation of the Hessian Matrix:

$$\mathbf{I}(\boldsymbol{\theta}) = -\mathbb{E}[\mathbf{H}(\mathbf{w}, \boldsymbol{\theta})]$$

- Why is it useful?
 - ▶ It can be used to assess whether the likelihood function is “well behaved” (Identification)
 - ▶ Important result: the information matrix is the inverse of the variance of the maximum likelihood estimator.
 - ▶ Cramér Rao lower bound.

Information matrix equality

Information matrix equality

The information matrix can be derived in two ways, either as minus the expected Hessian, or alternative as the variance of the score function, both evaluated at the true parameter θ_0

Information Identity

Assumption V: Finite Information

$\text{Var} \left[\frac{\partial}{\partial \boldsymbol{\theta}} \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \right] \equiv \text{Var} [\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})]$ exists.

Lemma (Information Identity)

Under Distribution, Differentiability and Finite Information Assumption:

$$\mathbb{E} \left[\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \ln f(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \right] = - \text{Var} [\mathbf{s}(\mathbf{w}; \boldsymbol{\theta})]$$

Proof: (Homework)

Information Identity

Note the following:

$$\begin{aligned}\text{Var} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] &= \mathbb{E} \left[\underbrace{\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)}_{(K \times 1)} \underbrace{\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)^\top}_{(1 \times K)} \right] + \underbrace{\mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] \mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^\top}_{=0} \\ &= \mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)^\top]\end{aligned}$$

Therefore we can write:

$$-\mathbf{I}(\boldsymbol{\theta}_0) = \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = -\text{Var} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = -\mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)^\top]$$

Example

Recall that:

$$\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = -0.5 \log(2\pi\sigma^2) - \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{2\sigma^2}$$

We have:

$$\begin{aligned} \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) &= \begin{pmatrix} \frac{1}{\sigma^2} \mathbf{x}_i \cdot \hat{\epsilon}_i \\ -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} \hat{\epsilon}_i^2 \end{pmatrix} \\ \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) &= \begin{pmatrix} -\frac{1}{\sigma^2} \mathbf{x}_i \mathbf{x}_i^\top & -\frac{1}{\sigma^4} \mathbf{x}_i \cdot \hat{\epsilon}_i \\ -\frac{1}{\sigma^4} \mathbf{x}_i^\top \cdot \hat{\epsilon}_i & \frac{1}{2\sigma^4} - \frac{1}{\sigma^6} \hat{\epsilon}_i^2 \end{pmatrix} \\ \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})^\top &= \begin{pmatrix} \frac{1}{\sigma^4} \mathbf{x}_i \mathbf{x}_i^\top \hat{\epsilon}_i^2 & -\frac{1}{2\sigma^4} \mathbf{x}_i \cdot \hat{\epsilon}_i + \frac{1}{2\sigma^6} \mathbf{x}_i \cdot \hat{\epsilon}_i^3 \\ -\frac{1}{2\sigma^4} \mathbf{x}_i^\top \cdot \hat{\epsilon}_i + \frac{1}{2\sigma^6} \mathbf{x}_i^\top \cdot \hat{\epsilon}_i^3 & \frac{1}{4\sigma^4} - \frac{1}{2\sigma^6} \hat{\epsilon}_i^2 + \frac{1}{4\sigma^8} \hat{\epsilon}_i^4 \end{pmatrix} \end{aligned}$$

where $\mathbf{w}_i = (y_i, \mathbf{x}_i^\top)^\top$, $\boldsymbol{\theta} = (\boldsymbol{\theta}^\top, \sigma^2)^\top$ and $\hat{\epsilon}_i \equiv y_i - \mathbf{x}_i^\top \boldsymbol{\beta}$

Example

So for $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ the $\hat{\epsilon}_i$ in these expressions can be replaced by ϵ_i . In the linear regression model, $\mathbb{E}(\epsilon_i | \mathbf{x}_i) = 0$. Also, since $\epsilon_i \sim N(0, \sigma_0^2)$, we have $\mathbb{E}(\epsilon_i^3) = 0$ and $\mathbb{E}(\epsilon_i^4) = 3\sigma_0^4$. Using these relations, we have:

$$-\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbb{E} [\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}) \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta})^\top] = \begin{pmatrix} \frac{1}{2\sigma_0^2} \mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top) & \mathbf{0} \\ \mathbf{0}^\top & \frac{1}{2\sigma_0^4} \end{pmatrix}$$

If $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top)$ is nonsingular, then $\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$ is nonsingular.

- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix

- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality

- 3 Estimation of Variance

- 4 Testing
 - Intuition
 - The Trinity
 - Proof, Proof and More Proof

Some Ideas

- For OLS estimator consistency can be shown by finding the sampling error function and applying LLN.
- This cannot be done for nonlinear estimator such as MLE since closed form solution for finite sample estimators do not exist.
- That is, the MLE is an implicit function of the random sample.

Question

How can we proceed?

Some Ideas

Using some LLN we know that:

$$\frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})] \quad (6)$$

for any fixed parameter value $\boldsymbol{\theta}$. That is, the sample average log-likelihood function converges to the expected log-likelihood for any value of $\boldsymbol{\theta}$. Recall that:

$$\begin{aligned} \hat{\boldsymbol{\theta}}_n &\equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \\ \boldsymbol{\theta}_0 &\equiv \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E} [\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})] \end{aligned}$$

We would like to say that, given that

$$\frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})], \text{ then } \hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_0$$

Some Ideas

We might be able to do this using the **continuous mapping theorem**.

- Let $X_n = \frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$,
- and $g(\cdot) = \arg \max_{\boldsymbol{\theta} \in \Theta}(\cdot)$

Then we would like to say that if $X_n \xrightarrow{P} X$ then $g(X_n) \xrightarrow{P} g_0(X)$.

In words:

If the sample average of the log likelihood function is close to the true expected value of the log likelihood function, then we would expect that $\hat{\boldsymbol{\theta}}_n$ will be close to the maximum of the expected likelihood (as n increases without bound)

However, we cannot do that!

What is the problem?

- The problem is that the argument of the $\arg \max_{\theta \in \Theta}(\cdot)$ is a function of θ , not a real vector:
 - ▶ The concept of convergence in probability was defined for **sequence of random variables**
- Therefore, we need to define what we mean by the probability limit of **sequence of random functions**, as opposed to a sequence of random variables:

Convergence for sequence of random variables $\implies X_n = X_n(\omega), \omega \in \Omega$

Convergence for sequence of random function $\implies Q_n = Q_n(\omega, \theta), \omega \in \Omega$

Example

Example

In ML estimation, the log-likelihood is a function of the sample data (a random vector that depends on ω) and of a parameter θ . By increasing the sample size, we obtain a sequence of log-likelihoods that depend on ω and θ .

Consistency

How is the distance between two functions over a set containing an infinite number of possible comparisons at different values of θ measured?

- IOW, since we are dealing with convergence on a **function space** we need to define when two functions are close to one another.
- To reduce the infinite dimensional character of a function to a one-dimensional concept of convergence, we take the supremum of the absolute difference of the function values over all θ in Θ

Uniform Convergence in Probability

Definition (Uniform Convergence in Probability)

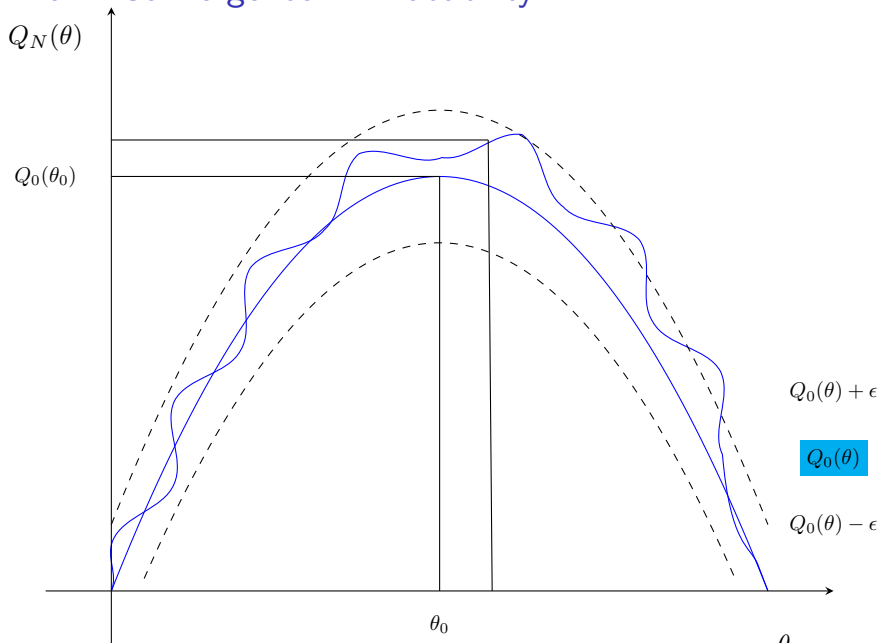
The sequence of real-valued functions $\{Q_n(\boldsymbol{\theta})\}$ converges uniformly in probability to the limit function $Q_0(\boldsymbol{\theta})$ if $\sup_{\boldsymbol{\theta} \in \Theta} |Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})| \xrightarrow{p} 0$. We will say that $Q_n(\boldsymbol{\theta}) \xrightarrow{p} Q_0(\boldsymbol{\theta})$ **uniformly**.

Another way to express uniform convergence in probability is:

$$\sup_{\boldsymbol{\theta} \in \Theta} |Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})| = o_p(1)$$

IOW, instead of requiring that the distance $|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})|$ converge in probability to 0 for each $\boldsymbol{\theta}$, we require convergence of $\sup_{\boldsymbol{\theta} \in \Theta} |Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})|$, which is the maximum distance that can be found by ranging over the space parameters.

Uniform Convergence in Probability



Uniform Convergence in Probability

Extending the concept to random vectors is straightforward. Now suppose that $\{Q_n(\boldsymbol{\theta})\}$ is a sequence of $K \times 1$ random vectors that depend both on the data and on the parameter $\boldsymbol{\theta} \in \Theta$. This sequence of **random vectors** is uniformly convergent in probability to $Q_0(\boldsymbol{\theta})$ if and only if

$$\sup_{\boldsymbol{\theta} \in \Theta} \|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})\| = o_p(1)$$

where $\|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})\|$ denotes the Euclidean norm of the vector $Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})$. By taking the supremum over $\boldsymbol{\theta}$ we obtain another random quantity that does not depend on $\boldsymbol{\theta}$.

Pointwise Convergence in probability

Definition (Pointwise Convergence in probability)

The sequence of real-valued functions $\{Q_n(\boldsymbol{\theta})\}$ converges pointwise in probability if and only if $|Q_n(\boldsymbol{\theta}) - Q_0(\boldsymbol{\theta})| \xrightarrow{p} 0$ **for each** $\boldsymbol{\theta} \in \Theta$

Uniform convergence is stronger than pointwise convergence.

Uniform LLN

Now we present the **uniform LLN** to study sequences of random functions which is analogous to the Chebychev's LLN for averages of random variables.

Theorem (Uniform LLN)

Suppose that $Q(\boldsymbol{\theta}, U)$ is continuous function over $\boldsymbol{\theta} \in \Theta$, a closed and bounded subset of \mathbb{R}^p , and that $\{U_n\}$ is a sequence of i.i.d. random variables with cdf $F_U(u)$. If $\mathbb{E}[\sup_{\boldsymbol{\theta} \in \Theta} \|Q(\boldsymbol{\theta}; U)\|]$ exists, then

- 1 $\mathbb{E}[Q(\boldsymbol{\theta}; U)]$ is continuous over $\boldsymbol{\theta} \in \Theta$ and,
- 2 $\frac{1}{n} \sum_{i=1}^n Q(\boldsymbol{\theta}; u_i) \xrightarrow{p} \mathbb{E}[Q(\boldsymbol{\theta}; U)]$ uniformly.

Or as Newey and McFadden (1994) (Lemma 2.4 state):

$$\sup_{\boldsymbol{\theta} \in \Theta} \left\| \frac{1}{n} \sum_{i=1}^n Q(\boldsymbol{\theta}; u_i) - \mathbb{E}[Q(\boldsymbol{\theta}; U)] \right\| \xrightarrow{p} \mathbf{0}$$

Uniform LLN

The following Theorem makes the connection between the uniform convergence of $\frac{1}{n} \sum_{i=1}^n Q(\boldsymbol{\theta}; u_i)$ to $\mathbb{E}[Q(\boldsymbol{\theta}; U)]$ and the convergence of $\hat{\boldsymbol{\theta}}_n$ to $\boldsymbol{\theta}_0$ using the assumption of **compact parameter space**.

Consistency

Theorem (Consistency of Maxima with Compact Parameter Space)

Suppose that:

- 1 (compact parameter space) $\Theta \subset \mathbb{R}^p$ is a closed and bounded parameter space,
- 2 (uniform convergence) $Q_n(\boldsymbol{\theta})$ is a sequence of function that convergence in probability uniformly to a function $Q_0(\boldsymbol{\theta})$ on Θ ,
- 3 (continuity) $Q_n(\boldsymbol{\theta})$ is continuous in $\boldsymbol{\theta}$ for any data $(\mathbf{w}_1, \dots, \mathbf{w}_n)$,
- 4 (identification) $Q_0(\boldsymbol{\theta})$ is uniquely maximized at $\boldsymbol{\theta}_0 \in \Theta$

then $\hat{\boldsymbol{\theta}}_n \equiv \arg \max_{\boldsymbol{\theta} \in \Theta} Q_n(\boldsymbol{\theta})$ converges in probability to $\boldsymbol{\theta}_0$.

Consistency

Intuition:

If $Q_n(\boldsymbol{\theta})$ converges uniformly to $Q_0(\boldsymbol{\theta})$, then the characteristics of $Q_n(\boldsymbol{\theta})$ will be close to the characteristics of $Q_0(\boldsymbol{\theta})$ as $n \rightarrow \infty$. One particular characteristic is the point $\boldsymbol{\theta}_0$ where $Q_0(\boldsymbol{\theta})$ is uniquely maximized. Then, it is expected that the maximizer of $Q_n(\boldsymbol{\theta})$, $\hat{\boldsymbol{\theta}}$, will be close to the maximizer of $Q_0(\boldsymbol{\theta})$.

Consistency

Theorem (Consistency of Maxima without Compactness)

Suppose that:

- 1 (interior) θ_0 is an element of the interior of a convex parameter space Θ ,
- 2 (pointwise convergence) $Q_n(\theta)$ converges in probability to $Q_0(\theta)$ for all $\theta \in \Theta$,
- 3 (concavity) $Q_n(\theta)$ is concave over the parameter space for any data $(\mathbf{w}_1, \dots, \mathbf{w}_n)$,
- 4 (identification) $Q_0(\theta)$ is uniquely maximized at $\theta_0 \in \Theta$

then, as $n \rightarrow \infty$, $\hat{\theta}_n$ exists with probability approaching 1 and $\hat{\theta}_n \xrightarrow{P} \theta_0$.

Consistency

Theorem (Consistency of conditional ML with compact parameter)

Let $\{y_i, \mathbf{x}_i\}$ be i.i.d with conditional density $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$ and let $\hat{\boldsymbol{\theta}}$ be the conditional ML estimator, which maximizes the average log conditional likelihood:

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$$

Suppose the model is correctly specified so that $\boldsymbol{\theta}_0$ is in Θ . Suppose that

- 1 (Compactness) the parameter space Θ is compact subset of \mathbb{R}^K ,
- 2 $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ is continuous in $\boldsymbol{\theta}$ for all (y_i, \mathbf{x}_i) (Here, note the Weierstrass's theorem),
- 3 $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ is measurable in (y_i, \mathbf{x}_i) for all $\boldsymbol{\theta} \in \Theta$ (so $\hat{\boldsymbol{\theta}}$ is well-defined random variable),
- 4 (identification) $\Pr [f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \neq f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)] > 0$ for all $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ in Θ ,
- 5 (dominance) $\mathbb{E} [\sup_{\boldsymbol{\theta} \in \Theta} |\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})|] < \infty$ (note: the expectation is over y_i and \mathbf{x}_i)

Then $\hat{\boldsymbol{\theta}} \xrightarrow{P} \boldsymbol{\theta}_0$

Sketch of Proof.

We would like to apply **Consistency of Maxima with Compact Parameter Space** Theorem. In this case, let $Q(\boldsymbol{\theta}; U) = \log f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$. Now we verify that the condition of the theorem are met:

- $f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ is continuous,
- Compactness states that Θ is a closed and bounded subset of E^K ,
- (y_i, \mathbf{x}_i) are i.i.d with conditional density $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$,
- Dominance I states that $\mathbb{E} \left[\sup_{\boldsymbol{\theta} \in \Theta} |\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})| \right]$ exists.

Therefore, $\mathbb{E} [\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$ is continuous and

$$\frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})] \quad (7)$$

uniformly. Let $Q_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})$ and $Q_0(\boldsymbol{\theta}) = \mathbb{E} [\log f(y_i|\mathbf{x}_i; \boldsymbol{\theta})]$. Under the additional assumption of Likelihood Identification, we can invoke the **strict expected log-likelihood inequality**: $\boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ that $\mathbb{E} [\log f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})] < \mathbb{E} [\log f(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}_0)]$. This implies that $Q_0(\boldsymbol{\theta})$ is uniquely maximized at $\boldsymbol{\theta}_0$. Therefore

$$\hat{\boldsymbol{\theta}}_n = \arg \max_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \xrightarrow{P} \boldsymbol{\theta}_0$$



Consistency

Theorem (Consistency of conditional ML without Compactness)

Let $\{y_i, \mathbf{x}_i\}$ be i.i.d with conditional density $f(y_i|\mathbf{x}_i; \theta_0)$ and let $\hat{\theta}$ be the conditional ML estimator, which maximizes the average log conditional likelihood:

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(y_i|\mathbf{x}_i; \theta)$$

Suppose the model is correctly specified so that θ_0 is in Θ . Suppose that

- 1 the true parameter vector θ_0 is an element of the interior of convex parameter space $\Theta \subset \mathbb{R}^K$,
- 2 $\log f(y_i|\mathbf{x}_i; \theta)$ is concave in θ for all (y_i, \mathbf{x}_i) ,
- 3 $\log f(y_i|\mathbf{x}_i; \theta)$ is measurable in (y_i, \mathbf{x}_i) ,
- 4 (identification) $\Pr [f(y_i|\mathbf{x}_i; \theta) \neq f(y_i|\mathbf{x}_i; \theta_0)] > 0$ for all $\theta \neq \theta_0$ in Θ ,
- 5 $\mathbb{E} [|\log f(y_i|\mathbf{x}_i; \theta)|] < \infty$ (i.e., $\mathbb{E} [\log f(y_i|\mathbf{x}_i; \theta)]$ exists and is finite) for all $\theta \in \Theta$ (note: the expectation is over y_i and \mathbf{x}_i)

Then, $n \rightarrow \infty$, $\hat{\theta}$ exists with probability approaching 1 and $\hat{\theta}_n \xrightarrow{p} \theta_0$

- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix

- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality

- 3 Estimation of Variance

- 4 Testing
 - Intuition
 - The Trinity
 - Proof, Proof and More Proof

Asymptotic

Theorem (Asymptotic Normality of Conditional ML)

Let $\mathbf{w} \equiv (y_i, \mathbf{x}_i^\top)^\top$ be iid. Suppose the conditions of either Theorem 18 or 19 are satisfied, so that $\hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_0$. Suppose, in addition, that:

- 1 $\boldsymbol{\theta}_0$ is in interior of Θ ,
- 2 $f(y_i|\mathbf{x}_i; \boldsymbol{\theta}_0)$ is twice continuously differentiable in $\boldsymbol{\theta}$ for all (y_i, \mathbf{x}_i) ,
- 3 $\mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbf{0}$ and $-\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)^\top]$,
- 4 (local dominance condition on the Hessian) for some neighborhood \mathcal{N} of $\boldsymbol{\theta}_0$,

$$\mathbb{E} \left[\sup_{\boldsymbol{\theta} \in \mathcal{N}} \|\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta})\| \right] < \infty$$

so that for any consistent estimator $\tilde{\boldsymbol{\theta}}$, $\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \tilde{\boldsymbol{\theta}}) \xrightarrow{P} \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$

- 5 $\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$ is nonsingular.

Then:

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N(\mathbf{0}, \mathbf{V}), \mathbf{V} = -\{\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]\}^{-1} = \{\mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)^\top]\}^{-1}$$

Asymptotic

The intuition is the following:

- Since usually we don't have an explicit solution for the estimator, we need to focus on the asymptotic behaviour of the score function.
- Assuming that $\hat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}$, the behaviour of the score function matters only within an arbitrary small neighbourhood of $\boldsymbol{\theta}_0$.
- ... after all, $\hat{\boldsymbol{\theta}}_n$ will fall within such neighborhoods with arbitrary high probability for a large enough sample size...
- ... and within such neighborhood the score function is essentially linear.. (Taylor series expansion)

Mean Value Theorem

Theorem (Mean Value Theorem)

Let $s : \mathbb{R}^K \rightarrow \mathbb{R}$ be defined on an open convex set $\Theta \subset \mathbb{R}^K$ such that s is continuously differentiable on Θ with gradient ∇s . Then for any points θ and θ_0 such that $s(\theta) = s(\theta_0) + \nabla s(\bar{\theta})(\theta - \theta_0)$

Asymptotic

Proof.

The objective function is:

$$Q_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i; \boldsymbol{\theta})$$

Given that $f(y_i | \mathbf{x}_i; \boldsymbol{\theta}_0)$ is twice continuously differentiable in $\boldsymbol{\theta}$, and given that $\boldsymbol{\theta}_0$ is in interior of Θ , then the maximum likelihood estimator satisfies

$$\frac{\partial \log L(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} = \mathbf{s}(\mathbf{w}; \hat{\boldsymbol{\theta}}) = \mathbf{0}$$

We need to know about the behavior of the gradient around the true parameter. Expand this set of equations in a Taylor series around the true parameters $\boldsymbol{\theta}_0$. We will use the mean value theorem to truncate the Taylor series at the second term,

$$\begin{aligned} \underbrace{\frac{\partial \log L(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}}}_{(K \times 1)} &= \underbrace{\frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}}}_{(K \times 1)} + \underbrace{\frac{\partial \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}}_{(K \times K)} \underbrace{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)}_{(K \times 1)} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) + \left[\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}) \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \end{aligned}$$

where $\bar{\boldsymbol{\theta}} = \alpha \hat{\boldsymbol{\theta}} + (1 - \alpha) \boldsymbol{\theta}_0$ for some $\alpha \in (0, 1)$



Proof.

So,

$$\sqrt{n} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \left[-\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}) \right]^{-1} \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \right)$$

We know that

$$\hat{\boldsymbol{\theta}} \xrightarrow{P} \boldsymbol{\theta}_0 \implies \bar{\boldsymbol{\theta}} \xrightarrow{P} \boldsymbol{\theta}_0$$

By uniform LLN, we know that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}) \xrightarrow{P} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta})] \quad \text{uniformly in } \boldsymbol{\theta} \in \Theta$$

Then, applying our Lemma:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}_n) \xrightarrow{P} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$$

since $\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$ exists. Finally, using **probability limit continuity** and **nonsingular information**, then:

$$\left[\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}_n) \right]^{-1} \xrightarrow{P} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1}$$



Proof.

Since (y_i, \mathbf{x}_i) are i.i.d. $\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\right)$ is the sum of variables $\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)$. The score identity lemma implies that $\mathbb{E}(\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)) = \mathbf{0}$, and the Information Identity implies that

$$\text{Var}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)] = \mathbb{E}[\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)\mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0)^\top] = -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]$$

The Lindberg-Levy CLT therefore implies:

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \right) \xrightarrow{d} \mathbf{N}(\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)])$$



Proof.

Then:

$$\begin{aligned}\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) &= \left[-\frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{w}_i; \bar{\boldsymbol{\theta}}) \right]^{-1} \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i; \boldsymbol{\theta}_0) \right) \\ &\xrightarrow{d} -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{N}(\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]) \\ &= \mathbf{N} \left[\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \right] \\ &= \mathbf{N} \left[\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \right] \\ &= \mathbf{N} \left[\mathbf{0}, [\mathbf{I}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \right]\end{aligned}$$



Variance Estimation

For large but finite samples, we can therefore write the approximate distribution of $\hat{\boldsymbol{\theta}}_n$ as

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

we have three potential estimators of $\mathbf{I}(\boldsymbol{\theta}_0)$:

- The empirical mean of minus the Hessian,

$$\hat{\mathbf{V}}^1 = \left(\frac{1}{n} \sum_{i=1}^n -\mathbf{H}(\mathbf{w}_i, \hat{\boldsymbol{\theta}}) \right)^{-1}$$

- The empirical variance of the score:

$$\hat{\mathbf{V}}^2 = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i, \hat{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{w}_i; \hat{\boldsymbol{\theta}})^\top \right)^{-1}$$

- Minus the expected Hessian evaluated at $\hat{\boldsymbol{\theta}}$:

$$\hat{\mathbf{V}}^3 = \left(\frac{1}{n} \sum_{i=1}^n -\mathbb{E} \left[\mathbf{H}(\mathbf{w}_i, \hat{\boldsymbol{\theta}}) \right] \right)^{-1}$$

Proof of Consistency

- Evaluated at a $\theta \in \Theta$, each estimator converges in probability uniformly to its expectation.
- Because $\hat{\theta}_n \xrightarrow{p} \theta_0$, evaluated at $\hat{\theta}_n$ each estimator converges in probability to $\mathbf{I}(\theta_0)$.
- Because matrix inversion is a continuous transformation, the inverse of each matrix is also a consistent estimator for the variance matrix of the asymptotic distribution of $\sqrt{n}(\hat{\theta}_n - \theta_0)$

- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix

- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality

- 3 Estimation of Variance

- 4 **Testing**
 - **Intuition**
 - The Trinity
 - Proof, Proof and More Proof

Hypothesis Testing

- ML estimator are distributed asymptotically normally:
 - ▶ As the sample size increases, the sampling distribution of an ML estimator becomes approximately normal

$$\hat{\beta} \stackrel{a}{\sim} N(\beta, \mathbf{V}_{\hat{\beta}})$$

where, for three coefficients:

$$\mathbf{V}_{\hat{\beta}} = \text{Var} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \sigma_{\hat{\beta}_0}^2 & \sigma_{\hat{\beta}_0, \hat{\beta}_1} & \sigma_{\hat{\beta}_0, \hat{\beta}_2} \\ \sigma_{\hat{\beta}_1, \hat{\beta}_0} & \sigma_{\hat{\beta}_1}^2 & \sigma_{\hat{\beta}_1, \hat{\beta}_2} \\ \sigma_{\hat{\beta}_2, \hat{\beta}_0} & \sigma_{\hat{\beta}_2, \hat{\beta}_1} & \sigma_{\hat{\beta}_2}^2 \end{pmatrix}$$

Hypothesis Testing

Consider the simple hypothesis

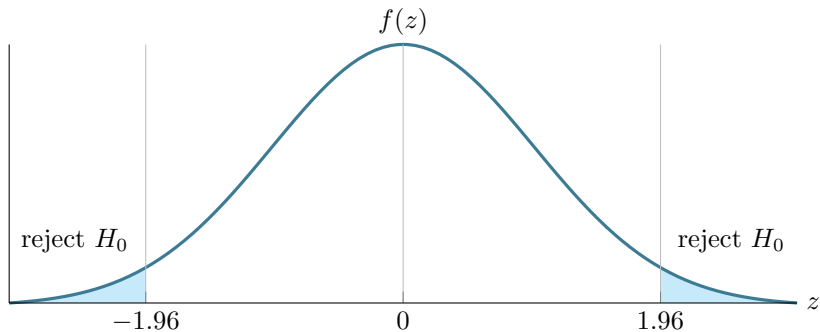
$$H_0 : \beta_k = \beta^*$$

where β^* is the hypothesized value, often equal to 0. Since $\sigma_{\hat{\beta}_k}^2$ is unknown, it must be estimated, which results in the test:

$$z = \frac{\hat{\beta}_k - \beta^*}{\hat{\sigma}_{\hat{\beta}_k}^2}$$

Under the assumptions justifying ML, if H_0 is true, then z is distributed approximately normally with mean of 0 and variance of 1 for large samples.

Testing



- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix

- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality

- 3 Estimation of Variance

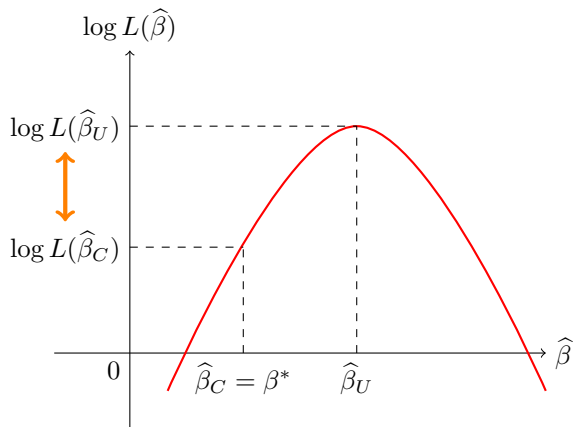
- 4 **Testing**
 - Intuition
 - **The Trinity**
 - Proof, Proof and More Proof

The Trinity

For more complex hypotheses we can use the Wald, likelihood ratio (LR), or Lagrange multiplier (LM) test. These test can be though of as a comparison between the estimates obtained after the constrains implied by the hypothesis have been imposed to the estimates obtained without the constraints.

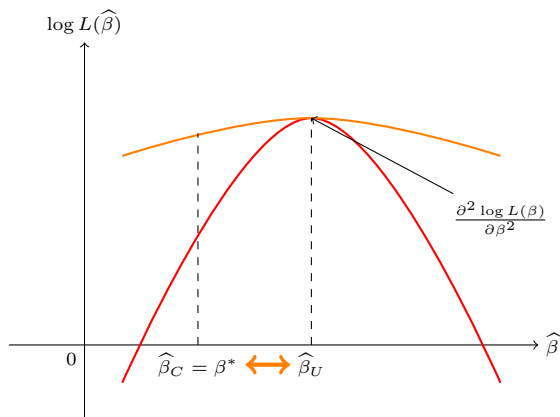
The Trinity: LR Test

- Log-likelihood function is the red solid line;
- $\hat{\beta}_U$: unconstrained estimator.
- The $H_0 : \beta = \beta^*$ imposes the constraint $\beta = \beta^*$, so that the constrained estimate is $\hat{\beta}_C = \beta^*$.
- Unless $\hat{\beta}_U = \beta^*$, $\ln L(\hat{\beta}_C) \leq \ln L(\hat{\beta}_U)$.
- If the constraint **significantly** reduces the likelihood, then the null hypothesis is rejected



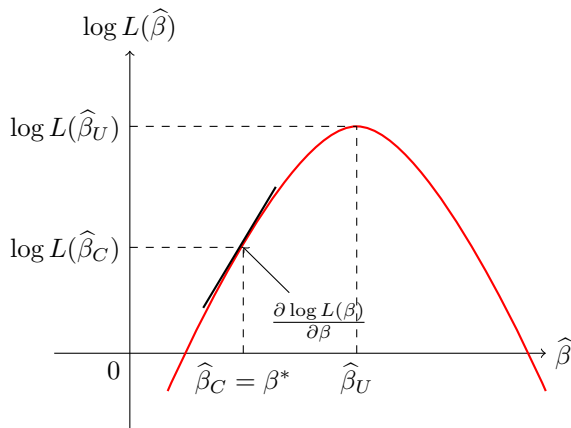
The Trinity: Wald Test

- The Wald test estimate the model **without constraints**, and assesses the constraint by considering 2 things:
 1. It measures the distance $\beta_U - \beta_C = \hat{\beta}_U - \hat{\beta}_C$.
 2. The distance is weighted by the curvature of the log likelihood function $\partial^2 \log L(\beta) / \partial \beta^2$
- The larger the distance, the less likely it is that the constraint is true.
- The larger the second derivative, the faster the curve is changing.
- The LL function (dashed line) is nearly flat, so the second derivative evaluated at $\hat{\beta}_U$ is relatively small.
- When the second derivative is small, the distance $\hat{\beta}_U$ and $\hat{\beta}_C$ is minor relative to the sampling variation.
- The second LL function has a larger second derivative.
- With a larger second derivative, the same distance between $\hat{\beta}_U$ and $\hat{\beta}_C$ might be significant.



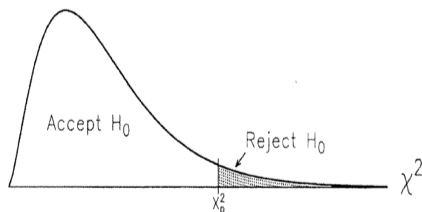
The Trinity: LM (Score) Test

- It only estimate the constrained model.
- It assesses the slope of the log likelihood function at the constraint.
- If H_0 is true, the slope (score) at the constraint should be close to 0.
- As with the Wald test, the curvature of the log likelihood function at the constraint is used to assess the significance of a nonzero slope.



The Trinity

- When H_0 is true, the Wald, LR, and LM tests are asymptotically equivalent.
- As $n \rightarrow \infty$, the sampling distributions of the three test $\xrightarrow{d} \chi_r^2$, where r is the number of constraints being tested.



They are similar only when $n \rightarrow \infty$. In small samples this is not necessary true.

Test Statistics

Assume that we have r ($r < K$) nonlinear restrictions (which includes linear restriction as special case)

$$\underbrace{\mathbf{r}(\boldsymbol{\theta}_0)}_{r \times 1} = \underbrace{\mathbf{0}}_{r \times 1}$$

The hypotheses are:

$$H_0 : \mathbf{r}(\boldsymbol{\theta}_0) = \mathbf{0}$$

$$H_1 : \mathbf{r}(\boldsymbol{\theta}_0) \neq \mathbf{0}$$

Also, denote $\hat{\boldsymbol{\theta}}$ as the unconstrained estimator and $\tilde{\boldsymbol{\theta}}$ as constrained estimator. So:

$$\tilde{\boldsymbol{\theta}}_{ML} = \arg \max_{\boldsymbol{\theta} \in \Theta} \left\{ \frac{1}{n} \sum_{i=1}^m \log f(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \right\} \quad s.t. \quad \mathbf{r}(\boldsymbol{\theta}_0) = \mathbf{0}$$

Test Statistics

Then, Wald, LM and LR statistics are defined as:

$$W = \underbrace{\mathbf{r}(\hat{\boldsymbol{\theta}})^\top}_{1 \times r} \left[\underbrace{\mathbf{R}(\hat{\boldsymbol{\theta}})}_{r \times K} \underbrace{\left(\frac{1}{n} \cdot \hat{\mathbf{V}} \right)}_{K \times K} \underbrace{\mathbf{R}(\hat{\boldsymbol{\theta}})^\top}_{K \times r} \right]^{-1} \underbrace{\mathbf{r}(\hat{\boldsymbol{\theta}})}_{r \times 1}$$
$$LM = \underbrace{\mathbf{s}(\tilde{\boldsymbol{\theta}})^\top}_{1 \times K} \underbrace{\left(\frac{1}{n} \cdot \tilde{\mathbf{V}} \right)}_{K \times K} \underbrace{\mathbf{s}(\tilde{\boldsymbol{\theta}})}_{K \times 1}$$
$$LR = 2 \cdot n \left(\underbrace{\frac{1}{n} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i; \hat{\boldsymbol{\theta}})}_{1 \times 1} - \underbrace{\frac{1}{n} \sum_{i=1}^n \ln f(y_i | \mathbf{x}_i; \tilde{\boldsymbol{\theta}})}_{1 \times 1} \right)$$

Test Statistics

where:

$$\underbrace{\mathbf{R}(\boldsymbol{\theta}_0)}_{r \times K} = \frac{\partial \mathbf{r}(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}^\top} \quad \text{Jacobian of } \mathbf{r}(\boldsymbol{\theta}_0)$$

$$\underbrace{\mathbf{V}}_{K \times K} = \mathbf{I}^{-1} = -\frac{1}{n} \sum_{i=1} \frac{\partial}{\partial \boldsymbol{\theta} \boldsymbol{\theta}^\top} \ln f(y_i | \mathbf{x}_i; \hat{\boldsymbol{\theta}})$$

- 1 Introduction
 - Motivation
 - Maximum Likelihood Estimator
 - Identification
 - The Score Function
 - The Information Matrix

- 2 Asymptotic Properties
 - Consistency
 - Asymptotic Normality

- 3 Estimation of Variance

- 4 **Testing**
 - Intuition
 - The Trinity
 - **Proof, Proof and More Proof**

Proof: Wald Statistic.

Write W as:

$$\mathbf{W} = \mathbf{c}_n^\top \mathbf{Z}_n^{-1} \mathbf{c}_n, \quad \mathbf{c}_n \equiv \sqrt{n} \mathbf{r}(\hat{\boldsymbol{\theta}}), \quad \mathbf{Z}_n \equiv \mathbf{R}(\hat{\boldsymbol{\theta}}) \widehat{\mathbf{V}} \mathbf{R}(\hat{\boldsymbol{\theta}})^\top \quad (8)$$

First, we will show that $\mathbf{c}_n \xrightarrow{d} N(\mathbf{0}, \mathbf{V}_f)$. Applying the MVT (21) (to truncate the Taylor series) to $\mathbf{r}(\hat{\boldsymbol{\theta}})$ around $\boldsymbol{\theta}_0$:

$$\mathbf{r}(\hat{\boldsymbol{\theta}}) = \underbrace{\mathbf{r}(\boldsymbol{\theta}_0)}_{=0 \text{ under } H_0} + \mathbf{R}(\bar{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$$

$$\begin{aligned} \underbrace{\sqrt{n} \mathbf{r}(\hat{\boldsymbol{\theta}})}_{r \times 1} &= \underbrace{\mathbf{R}(\bar{\boldsymbol{\theta}})}_{r \times K} \underbrace{\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)}_{K \times 1} \quad \text{multiplying by } \sqrt{n} \\ &= \mathbf{R}(\bar{\boldsymbol{\theta}}) \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \mathbf{R}(\boldsymbol{\theta}_0) \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) - \mathbf{R}(\boldsymbol{\theta}_0) \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ &= (\mathbf{R}(\bar{\boldsymbol{\theta}}) - \mathbf{R}(\boldsymbol{\theta}_0)) \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \mathbf{R}(\boldsymbol{\theta}_0) \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \end{aligned}$$

where $\bar{\boldsymbol{\theta}} = \alpha \hat{\boldsymbol{\theta}} + (1 - \alpha) \boldsymbol{\theta}_0$



Proof: Wald Statistic.

Note that

$$\because \hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0 \implies \bar{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$$

$$\because F(\cdot) \text{ is continuous} \implies F(\bar{\boldsymbol{\theta}}) \xrightarrow{p} F(\boldsymbol{\theta}_0)$$

Note that $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d}$ some random variable, then we can write:

$$\sqrt{n}\mathbf{r}(\hat{\boldsymbol{\theta}}) = \underbrace{(\mathbf{R}(\bar{\boldsymbol{\theta}}) - \mathbf{R}(\boldsymbol{\theta}_0))}_{\xrightarrow{p} \mathbf{0}} \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \mathbf{R}(\boldsymbol{\theta}_0)\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$$

$$\sqrt{n}\mathbf{r}(\hat{\boldsymbol{\theta}}) = \mathbf{R}(\boldsymbol{\theta}_0)\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + o_p(1)$$

Furthermore, we know that:

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i, \boldsymbol{\theta}_0) + o_p(1) \quad (9)$$



Proof: Wald Statistic.

Then:

$$\begin{aligned}\sqrt{nr}(\hat{\boldsymbol{\theta}}) &= \mathbf{R}(\boldsymbol{\theta}_0)\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + o_p(1) \\ &= \mathbf{R}(\boldsymbol{\theta}_0) \left[\mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i, \boldsymbol{\theta}_0) + o_p(1) \right] + o_p(1) \\ &= \mathbf{R}(\boldsymbol{\theta}_0) \mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]^{-1} \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{s}(\mathbf{w}_i, \boldsymbol{\theta}_0)}_{\xrightarrow{d} \mathbf{N}(\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)])} + \underbrace{o_p(1)}_{o_p(1)+o_p(1)}\end{aligned}$$

Then:

$$\begin{aligned}\sqrt{nr}(\hat{\boldsymbol{\theta}}) &\xrightarrow{d} \mathbf{N} \left(\mathbf{0}, -\mathbf{R}(\boldsymbol{\theta}_0) \mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]^{-1} \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] \mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]^{-1} \mathbf{R}(\boldsymbol{\theta}_0)^\top \right) \\ &\xrightarrow{d} \mathbf{N} \left(\mathbf{0}, \mathbf{R}(\boldsymbol{\theta}_0) \underbrace{\left[-\mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]^{-1} \right]}_{\mathbf{V}} \mathbf{R}(\boldsymbol{\theta}_0)' \right)\end{aligned}$$



Proof: Wald Statistic.

It follows that:

$$\mathbf{W} = \mathbf{c}_n^\top \mathbf{Z}_n^{-1} \mathbf{c}_n = \sqrt{n} \mathbf{r}(\hat{\boldsymbol{\theta}}) [\mathbf{R}(\boldsymbol{\theta}_0) \mathbf{V}_0 \mathbf{R}(\boldsymbol{\theta}_0)^\top]^{-1} \sqrt{n} \mathbf{r}(\hat{\boldsymbol{\theta}}) \xrightarrow{d} \chi^2(\#r)$$

If we have consistent estimators of $\mathbf{R}(\boldsymbol{\theta}_0)$ and \mathbf{V} , then limit results for continuous functions imply that:

$$\sqrt{n} \mathbf{r}(\hat{\boldsymbol{\theta}}) \left[\hat{\mathbf{R}}(\hat{\boldsymbol{\theta}}_0) \hat{\mathbf{V}} \hat{\mathbf{R}}(\hat{\boldsymbol{\theta}}_0)^\top \right]^{-1} \sqrt{n} \mathbf{r}(\hat{\boldsymbol{\theta}}) \xrightarrow{d} \chi^2(\#r)$$



Preliminaries to the next two statistics

The Lagrangian for the constraint is

$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i, \boldsymbol{\theta}) - \underbrace{\boldsymbol{\lambda}^\top}_{1 \times r} \underbrace{\mathbf{r}(\boldsymbol{\theta})}_{r \times 1}$$

Then FOC:

$$\begin{aligned} \sqrt{n} \mathbf{s}(\tilde{\boldsymbol{\theta}}) + \sqrt{n} \mathbf{R}(\tilde{\boldsymbol{\theta}})^\top \boldsymbol{\lambda} &= \mathbf{0} \\ \sqrt{n} \mathbf{r}(\tilde{\boldsymbol{\theta}}) &= \mathbf{0} \end{aligned} \tag{10}$$

Following a similar reasoning as in previous proof:

$$\sqrt{n} \mathbf{r}(\tilde{\boldsymbol{\theta}}) = \mathbf{R}(\boldsymbol{\theta}_0) \sqrt{n}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + o_p(1) = O_p(1)$$

Preliminaries to the next two statistics

A Taylor expansion of $\mathbf{s}(\tilde{\boldsymbol{\theta}})$ yields:

$$\begin{aligned}\frac{\partial \log L(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} &= \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + \frac{\partial^2 \log L(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ &= \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + \frac{\partial^2 \log L(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) - \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ &= \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \left(\frac{\partial^2 \log L(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} - \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right) (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ \sqrt{n} \frac{\partial \log L(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} &= \sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \\ &\quad + \underbrace{\left(\frac{\partial^2 \log L(\tilde{\boldsymbol{\theta}})}{n^{-1} \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} - n^{-1} \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right)}_{o_p(1)} \underbrace{\sqrt{n} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)}_{O_p(1)} \\ &= \underbrace{\sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}}}_{\xrightarrow{d} \mathbf{N}(\mathbf{0}, -\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)])} + \underbrace{\frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top}}_{\mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + o_p(1) \\ &= O_p(1) + O_p(1)O_p(1) + o_p(1) \\ &= O_p(1) + O_p(1) + o_p(1) = O_p(1)\end{aligned}$$

Preliminaries to the next two statistics

Then

$$\mathbf{R}(\tilde{\boldsymbol{\theta}})^\top \sqrt{n} \boldsymbol{\lambda}_n = - \underbrace{\sqrt{n} \mathbf{s}(\tilde{\boldsymbol{\theta}})}_{O_p(1)} \implies \mathbf{R}(\tilde{\boldsymbol{\theta}})' \sqrt{n} \boldsymbol{\lambda}_n = O_p(1)$$

Since $\mathbf{R}(\tilde{\boldsymbol{\theta}}) \xrightarrow{p} \mathbf{R}_0$, we obtain:

$$\mathbf{R}(\tilde{\boldsymbol{\theta}}) \sqrt{n} \boldsymbol{\lambda}_n = \mathbf{R}_0^\top \sqrt{n} \boldsymbol{\lambda}_n + \left(\mathbf{R}(\tilde{\boldsymbol{\theta}}) - \mathbf{R}_0 \right)' \sqrt{n} \boldsymbol{\lambda}_n = \mathbf{R}_0^\top \sqrt{n} \boldsymbol{\lambda}_n + o_p(1)$$

Substituting these three equations into the FOCs, we obtain:

$$\begin{pmatrix} \mathbb{E}[\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)] & \mathbf{R}_0^\top \\ (K \times K) & K \times r \\ \mathbf{R}_0 & \mathbf{0} \\ (r \times K) & (r \times r) \end{pmatrix} \begin{pmatrix} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ (K \times 1) \\ \sqrt{n} \boldsymbol{\lambda}_n \\ (r \times 1) \end{pmatrix} = \begin{pmatrix} -\sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \\ (K \times 1) \\ \mathbf{0} \\ (r \times 1) \end{pmatrix} + o_p(1)$$

Preliminaries to the next two statistics

This can be solved using the partitioned inverse formula:

$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1} \mathbf{A}_{12} (\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12})^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1} \mathbf{A}_{12} (\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12})^{-1} \\ -(\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12})^{-1} \mathbf{A}_{21} \mathbf{A}_{11}^{-1} & (\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12})^{-1} \end{pmatrix}$$

Then:

$$\begin{pmatrix} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ (K \times 1) \\ \sqrt{n} \boldsymbol{\lambda}_n \\ (r \times 1) \end{pmatrix} = \begin{pmatrix} -\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} + \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top (\mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top)^{-1} \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \\ -(\mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top)^{-1} \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \end{pmatrix}$$

Proof LR.

By second order Taylor series:

$$\log L(\tilde{\boldsymbol{\theta}}) = \log L(\hat{\boldsymbol{\theta}}) + \frac{\partial \log L(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + \frac{1}{2} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \frac{\partial^2 \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$$

where $\bar{\boldsymbol{\theta}} = \alpha \tilde{\boldsymbol{\theta}} + (1 - \alpha) \hat{\boldsymbol{\theta}}$ for some $\alpha \in \alpha [0, 1]$. Recall that $\frac{\partial \log L(\hat{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} = \mathbf{0}$ and $n^{-1} \frac{\partial^2 \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \xrightarrow{p} \mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)]$. It follows that:

$$\log L(\tilde{\boldsymbol{\theta}}) - \log L(\hat{\boldsymbol{\theta}}) = \frac{1}{2} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \frac{\partial^2 \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$$

$$2 \cdot (\log L(\tilde{\boldsymbol{\theta}}) - \log L(\hat{\boldsymbol{\theta}})) = (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})^\top \frac{\partial^2 \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$$

$$2\sqrt{n}\sqrt{n} \cdot (\log L(\tilde{\boldsymbol{\theta}}) - \log L(\hat{\boldsymbol{\theta}})) = \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' n^{-1} \frac{\partial^2 \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$$

$$2 \cdot n \cdot (\log L(\tilde{\boldsymbol{\theta}}) - \log L(\hat{\boldsymbol{\theta}})) = \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})^\top n^{-1} \frac{\partial^2 \log L(\bar{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$$



Proof LR.

Adding and subtracting $\sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})^\top n^{-1} \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})$:

$$\begin{aligned} 2 \cdot n \cdot (\log L(\tilde{\boldsymbol{\theta}}) - \log L(\hat{\boldsymbol{\theta}})) &= \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})^\top n^{-1} \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + \\ &\quad + \underbrace{\sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})^\top \left[n^{-1} \frac{\partial^2 \log L(\tilde{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} - n^{-1} \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right]}_{o_p(1) = O_p(1) o_p(1) O_p(1)} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) \\ &= \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})^\top n^{-1} \frac{\partial^2 \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + o_p(1) \\ &= \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \mathbb{E}[\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)] \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + o_p(1) \\ 2 \cdot n \cdot (\log L(\hat{\boldsymbol{\theta}}) - \log L(\tilde{\boldsymbol{\theta}})) &= -\sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \mathbb{E}[\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)] \sqrt{n} (\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + o_p(1) \quad \times -1 \end{aligned}$$



Proof LR.

We know that

$$\begin{aligned}\sqrt{n}(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) &= \sqrt{n}(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) - \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ &= - \left(\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} - \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top (\mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top)^{-1} \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \right) \times \\ &\times \sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} - \left(-\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right) + o_p(1) \\ &= \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top (\mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top)^{-1} \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} + o_p(1)\end{aligned}$$

Then:

$$\begin{aligned}2 \cdot n \cdot (\log L(\hat{\boldsymbol{\theta}}) - \log L(\tilde{\boldsymbol{\theta}})) &= -\sqrt{n}(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}})' \mathbb{E} [\mathbf{H}(\mathbf{w}_i, \boldsymbol{\theta}_0)] \sqrt{n}(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}) + o_p(1) \\ &= \left(\sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right)' \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top \left(\mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \mathbf{R}_0^\top \right)^{-1} \times \\ &= \times \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)]^{-1} \underbrace{\left(\sqrt{n} \frac{\partial \log L(\boldsymbol{\theta}_0)}{\partial \boldsymbol{\theta}} \right)}_{\xrightarrow{d} \mathbf{N}(\mathbf{0}, -\mathbb{E} [\mathbf{H}(\mathbf{w}_i; \boldsymbol{\theta}_0)])}\end{aligned}$$



Proof LR.

Then :

$$\begin{aligned} \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \theta_0)]^{-1} \left(\sqrt{n} \frac{\partial \log L(\theta_0)}{\partial \theta} \right) &\xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \theta_0)]^{-1} - \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \theta_0)] \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \theta_0)]^{-1} \mathbf{R}_0^T \right) \\ &\xrightarrow{d} \mathcal{N} \left(\mathbf{0}, -\mathbf{R}_0 \mathbb{E} [\mathbf{H}(\mathbf{w}_i; \theta_0)]^{-1} \mathbf{R}_0' \right) \end{aligned}$$

This asymptotic variance cancels against the central term of the quadratic form, and hence we are looking at the norm of a $\#r$ -dimensional standard normal vector:

$$LR \equiv 2 \cdot n \cdot \left(\log L(\hat{\theta}) - \log L(\tilde{\theta}) \right) \xrightarrow{d} \chi^2(\#r)$$

