

1. Introduction

1.1 Why empirical likelihood

- nonparametric method: without having to assume the form of the underlying distribution
- likelihood based inference: taking the advantages of likelihood methods
- alternative method when other (more conventional) methods are not applicable

Remark. (i) For $N(\mu, \sigma^2)$, $\gamma = 0$ and $\kappa = 0$.

(ii) For symmetric distributions, $\gamma = 0$.

(iii) When $\kappa > 0$, heavier tails than those of $N(\mu, \sigma^2)$.

Example 1. Somites of earthworms.

Earthworms have segmented bodies. The segments are known as somites. As a worm grows, both the number and the length of its somites increases.

The dataset contains the No. of somites on each of 487 worms gathered near Ann Arbor in 1902.

The histogram shows that the distribution is *skewed to the left*, and has a *heavier tail to the left*.

Skewness: $\gamma = \frac{E\{(X-EX)^3\}}{\{\text{Var}(X)\}^{3/2}}$, — a measure for symmetry

Kurtosis: $\kappa = \frac{E\{(X-EX)^4\}}{\{\text{Var}(X)\}^2} - 3$, — a measure for tail-heaviness

Estimation for γ and κ

Let $\bar{X} = n^{-1} \sum_{1 \leq i \leq n} X_i$, and $\hat{\sigma}^2 = (n-1)^{-1} \sum_{1 \leq i \leq n} (X_i - \bar{X})^2$.

$$\hat{\gamma} = \frac{1}{n\hat{\sigma}^3} \sum_{i=1}^n (X_i - \bar{X})^3, \quad \hat{\kappa} = \frac{1}{n\hat{\sigma}^4} \sum_{i=1}^n (X_i - \bar{X})^4.$$

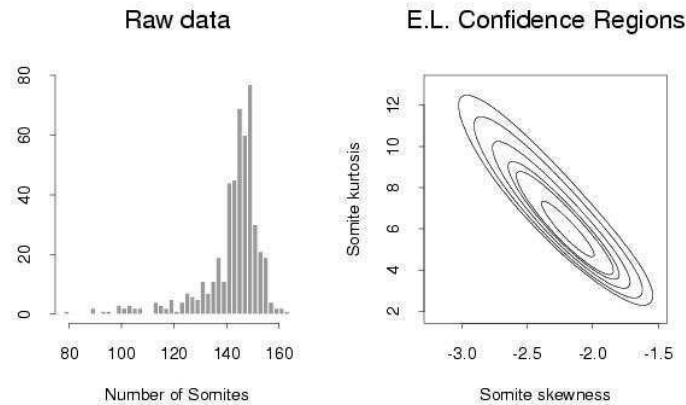
How to find the confidence sets for (γ, κ) ?

Answer: Empirical likelihood contours.

Let $l(\gamma, \kappa)$ be the (log-) empirical likelihood function of (γ, κ) . The confidence region for (γ, κ) is defined as

$$\{(\gamma, \kappa) : l(\gamma, \kappa) > C\},$$

where $C > 0$ is a constant determined by the confidence level, i.e. $P\{l(\gamma, \kappa) > C\} = 1 - \alpha$.



In the second panel, the empirical likelihood confidence regions (i.e. contours) correspond to confidence levels of 50%, 90%, 95%, 99%, 99.9% and 99.99%.

Note. $(\gamma, \kappa) = (0, 0)$ is not contained in the confidence regions

1.2 Introducing empirical likelihood

Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a random sample from an unknown distribution $F(\cdot)$. We *know nothing* about $F(\cdot)$.

In practice we observe $X_i = x_i$ ($i = 1, \dots, n$), x_1, \dots, x_n are n known numbers.

Basic idea. Assume F is a discrete distribution on $\{x_1, \dots, x_n\}$ with

$$p_i = F(x_i), \quad i = 1, \dots, n,$$

where

$$p_i \geq 0, \quad \sum_{i=1}^n p_i = 1.$$

What is the likelihood function of $\{p_i\}$? What is the MLE?

Why do conventional methods not apply?

Parametric likelihood. **Not normal distribution!** Likelihood inference for high moments is typically **not robust** wrt a misspecified distribution.

Bootstrap. Difficult in **picking out the confidence region** from a point cloud consisting of a large number of bootstrap estimates for (γ, κ) .

For example, given 1000 bootstrap estimates for (γ, κ) , ideally 95% confidence region should contain 950 central points.

In practice, we restrict to rectangle or ellipse regions in order to facilitate the estimation.

Since

$$P\{X_1 = x_1, \dots, X_n = x_n\} = p_1 \cdots p_n,$$

the likelihood is

$$L(p_1, \dots, p_n) \equiv L(p_1, \dots, p_n; \mathbf{X}) = \prod_{i=1}^n p_i,$$

which is called an *empirical likelihood*.

Remark. The number of parameters is the same as the number of observations.

Note

$$\left(\prod_{i=1}^n p_i \right)^{1/n} \leq \frac{1}{n} \sum_{i=1}^n p_i = \frac{1}{n},$$

the equality holds iff $p_1 = \dots = p_n = 1/n$.

Put $\hat{p}_i = 1/n$, we have

$$L(p_1, \dots, p_n; \mathbf{X}) \leq L(\hat{p}_1, \dots, \hat{p}_n; \mathbf{X})$$

for any $p_i \geq 0$ and $\sum_i p_i = 1$.

Hence the MLE based on the empirical likelihood, which is called **maximum empirical likelihood estimator (MELE)**, puts the equal probability mass $1/n$ on the n observed values x_1, \dots, x_n .

Namely the MELE for F is the uniform distribution on observed data points. The corresponding distribution function

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

is called the **empirical distribution** of the sample $\mathbf{X} = (X_1, \dots, X_n)^T$.

Remarks. (i) MELEs, without further constraints, are simply the method of moments estimators, which is not new.

(ii) Empirical likelihood is a powerful tool in dealing with testing hypotheses and interval estimation in [a nonparametric manner](#) based on the *likelihood* tradition, which also involves evaluating MELEs under some further constraints.

Example 2. Find the MELE for $\mu \equiv EX_1$.

Corresponding to the EL,

$$\mu = \sum_{i=1}^n p_i x_i = \mu(p_1, \dots, p_n).$$

Therefore, the MELE for μ is

$$\hat{\mu} = \mu(\hat{p}_1, \dots, \hat{p}_n) = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}.$$

Similarly, the MELE for $\mu_k \equiv E(X_1^k)$ is simply the sample k -th moment:

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k.$$

2. Empirical likelihood for means

Let X_1, \dots, X_n be a random sample from an unknown distribution.

Goal: test hypotheses on $\mu \equiv EX_1$, or find confidence intervals for μ .

Tool: *empirical likelihood ratios (ELR)*

2.1 Tests Consider the hypotheses

$$H_0: \mu = \mu_0 \quad \text{vs} \quad H_1: \mu \neq \mu_0.$$

Let $L(p_1, \dots, p_n) = \prod_i p_i$. We reject H_0 for large values of the ELR

$$T = \frac{\max L(p_1, \dots, p_n)}{\max_{H_0} L(p_1, \dots, p_n)} = \frac{L(n^{-1}, \dots, n^{-1})}{L(\tilde{p}_1, \dots, \tilde{p}_n)},$$

where $\{\tilde{p}_1\}$ are the **constrained MELEs** for $\{p_i\}$ under H_0 .

Two problems:

- (i) $\tilde{p}_i = ?$
- (ii) What is the distribution of T under H_0 ?

(i) The constrained MELEs $\tilde{p}_i = p_i(\mu_0)$, where $\{p_i(\mu)\}$ are the solution of the maximisation problem:

$$\max_{\{p_i\}} \sum_{i=1}^n \log p_i$$

subject to the conditions

$$p_i \geq 0, \quad \sum_{i=1}^n p_i = 1, \quad \sum_{i=1}^n p_i x_i = \mu.$$

The solution for the above problem is given in the theorem below. Note

$$x_{(1)} \equiv \min_i x_i \leq \sum_{i=1}^n p_i x_i \leq \max_i x_i \equiv x_{(n)}.$$

It is natural we require $x_{(1)} \leq \mu \leq x_{(n)}$.

have

$$p_i^{-1} + \psi + \lambda x_i = 0 \quad (3)$$

$$\sum_i p_i = 1 \quad (4)$$

$$\sum_i p_i x_i = \mu \quad (5)$$

By (3),

$$p_i = -1/(\psi + \lambda x_i). \quad (6)$$

Hence, $1 + \psi p_i + \lambda x_i p_i = 0$, which implies $\psi = -(n + \lambda \mu)$. This together with (6) imply (1). By (1) and (5),

$$\sum_i \frac{x_i}{n - \lambda(x_i - \mu)} = \mu. \quad (7)$$

It follows (4) that

$$\mu = \mu \sum_i p_i = \sum_i \frac{\mu}{n - \lambda(x_i - \mu)}.$$

This together with (7) imply (2).

Theorem 1. For $\mu \in (x_{(1)}, x_{(n)})$,

$$p_i(\mu) = \frac{1}{n - \lambda(x_i - \mu)} > 0, \quad 1 \leq i \leq n, \quad (1)$$

where λ is the unique solution of the equation

$$\sum_{j=1}^n \frac{x_j - \mu}{n - \lambda(x_j - \mu)} = 0 \quad (2)$$

in the interval $(\frac{n}{x_{(1)} - \mu}, \frac{n}{x_{(n)} - \mu})$.

Proof. We use the Lagrange multiplier technique to solve this optimisation problem. Put

$$Q = \sum_i \log p_i + \psi \left(\sum_i p_i - 1 \right) + \lambda \left(\sum_i p_i x_i - \mu \right).$$

Letting the partial derivatives of Q w.r.t. p_i , ψ and λ equal 0, we

Now let $g(\lambda)$ be the function on the LHS of (2). Then

$$g(\lambda) = \sum_i \frac{(x_i - \mu)^2}{\{n - \lambda(x_i - \mu)\}^2} > 0.$$

Hence $g(\lambda)$ is a strictly increasing function. Note

$$\lim_{\lambda \uparrow \frac{n}{x_{(n)} - \mu}} g(\lambda) = \infty, \quad \lim_{\lambda \downarrow \frac{n}{x_{(1)} - \mu}} g(\lambda) = -\infty,$$

Hence $g(\lambda) = 0$ has a unique solution between in the interval

$$\left(\frac{n}{x_{(1)} - \mu}, \frac{n}{x_{(n)} - \mu} \right).$$

Note for any λ in this interval,

$$\frac{1}{n - \lambda(x_{(1)} - \mu)} > 0, \quad \frac{1}{n - \lambda(x_{(n)} - \mu)} > 0,$$

and $1/\{n - \lambda(x - \mu)\}$ is a monotonic function of x . It holds that $p_i(\mu) > 0$ for all $1 \leq i \leq n$.

Remarks. (a) When $\mu = \bar{x} = \bar{X}$, $\lambda = 0$, and

$$p_i(\mu) = 1/n, \quad i = 1, \dots, n.$$

It may be shown for μ close to $E(X_i)$, and n large

$$p_i(\mu) \approx \frac{1}{n} \cdot \frac{1}{1 + \frac{\bar{x} - \mu}{S(\mu)}(x_i - \mu)},$$

where $S(\mu) = \frac{1}{n} \sum_i (x_i - \mu)^2$.

(b) We may view

$$L(\mu) = L\{p_1(\mu), \dots, p_n(\mu)\}.$$

as a **profile empirical likelihood** for μ .

Hypothetically consider an 1-1 parameter transformation from $\{p_1, \dots, p_n\}$ to $\{\mu, \theta_1, \dots, \theta_{n-1}\}$. Then

$$L(\mu) = \max_{\{\theta_i\}} L(\mu, \theta_1, \dots, \theta_{n-1}) = L\{\mu, \hat{\theta}_1(\mu), \dots, \hat{\theta}_{n-1}(\mu)\}$$

the ELR statistic is

$$\begin{aligned} T &= \frac{\max L(p_1, \dots, p_n)}{\max_{H_0} L(p_1, \dots, p_n)} = \frac{(1/n)^n}{L(\mu_0)} \\ &= \prod_{i=1}^n \frac{1}{np_i(\mu_0)} = \prod_{i=1}^n \left\{1 - \frac{\lambda}{n}(X_i - \mu_0)\right\}. \end{aligned}$$

where λ is the unique solution of

$$\sum_{j=1}^n \frac{X_j - \mu_0}{n - \lambda(X_j - \mu_0)} = 0.$$

Theorem 2. Let $E(X_1^2) < \infty$. Then under H_0 ,

$$2 \log T = 2 \sum_{i=1}^n \log \left\{1 - \frac{\lambda}{n}(X_i - \mu_0)\right\} \rightarrow \chi_1^2$$

in distribution as $n \rightarrow \infty$.

A sketch proof. Under H_0 , $EX_i = \mu_0$. Therefore μ_0 is close to \bar{X} for large n . Hence the λ , or more precisely, $\lambda_n \equiv \lambda/n$ is small,

(c) The likelihood function $L(\mu)$ may be calculated using R-code and Splus-code, downloaded at

<http://www-stat.stanford.edu/~owen/empirical/>

(ii) The asymptotic theorem for the classic likelihood ratio tests (i.e. Wilks' Theorem) still holds for the ELR tests.

Let X_1, \dots, X_n i.i.d., and $\mu = E(X_1)$. To test

$$H_0: \mu = \mu_0 \quad \text{vs} \quad H_1: \mu \neq \mu_0,$$

which is the solution of $f(\lambda_n) = 0$, where

$$f(\lambda_n) = \frac{1}{n} \sum_{j=1}^n \frac{X_j - \mu_0}{1 - \lambda_n(X_j - \mu_0)}.$$

By a simple Taylor expansion $0 = f(\lambda_n) \approx f(0) + f'(0)\lambda_n$,

$$\lambda_n \approx -f(0)/f'(0) = -(\bar{X} - \mu_0) / \left(\frac{1}{n} \sum_j (X_j - \mu_0)^2\right).$$

Now

$$\begin{aligned} 2 \log T &\approx 2 \sum_i \left\{-\lambda_n(X_i - \mu_0) - \frac{\lambda_n^2}{2}(X_i - \mu_0)^2\right\} \\ &= -2\lambda_n n(\bar{X} - \mu_0) - \lambda_n^2 \sum_i (X_i - \mu_0)^2 \approx \frac{n(\bar{X} - \mu_0)^2}{n^{-1} \sum_i (X_i - \mu_0)^2}. \end{aligned}$$

By the LLN, $n^{-1} \sum_i (X_i - \mu_0)^2 \rightarrow \text{Var}(X_1)$. By the CLT, $\sqrt{n}(\bar{X} - \mu_0) \rightarrow N(0, \text{Var}(X_1))$ in distribution. Hence $2 \log T \rightarrow \chi_1^2$ in distribution.

2.2 Confidence intervals for μ .

For a given $\alpha \in (0, 1)$, since we will not reject the null hypothesis

$$H_0 : \mu = \mu_0$$

iff $2 \log T < \chi_{1,1-\alpha}^2$, where $P\{\chi_1^2 \leq \chi_{1,1-\alpha}^2\} = 1 - \alpha$. For $\alpha = 0.05$, $\chi_{1,1-\alpha}^2 = 3.84$.

Hence a $100(1 - \alpha)\%$ confidence interval for μ is

$$\begin{aligned} & \left\{ \mu \mid -2 \log\{L(\mu)n^n\} < \chi_{1,1-\alpha}^2 \right\} \\ &= \left\{ \mu \mid \sum_{i=1}^n \log p_i(\mu) > -0.5\chi_{1,1-\alpha}^2 - n \log n \right\} \\ &= \left\{ \mu \mid \sum_{i=1}^n \log\{np_i(\mu)\} > -0.5\chi_{1,1-\alpha}^2 \right\}. \end{aligned}$$

Let $\mu = EX_i$.

$$H_0 : \mu = 0 \quad vs \quad H_1 : \mu > 0$$

(i) **Standard approach:** Assume $\{X_1, \dots, X_{15}\}$ is a random sample from $N(\mu, \sigma^2)$

MLE: $\hat{\mu} = \bar{X} = 2.61$

The t -test statistic:

$$T = \sqrt{n}\bar{X}/s = 2.14$$

Since $T \sim t(14)$ under H_0 , the p -value is 0.06 — *significant but not overwhelming*.

Is $N(\mu, \sigma^2)$ an appropriate assumption? as the data do not appear to be normal (with a heavy left tail); see Fig(a).

Example 3. Darwin's data: gains in height of plants from cross-fertilisation

$$X = \text{height}(\text{Cross-F}) - \text{height}(\text{Self-F})$$

15 observations:

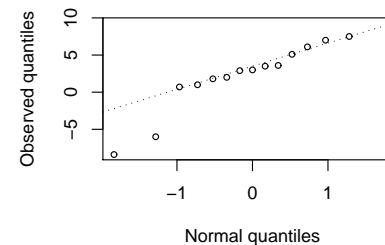
6.1, -8.4, 1.0, 2.0, 0.7, 2.9, 3.5, 5.1, 1.8, 3.6, 7.0, 3.0, 9.3, 7.5, -6.0

The sample mean $\bar{X} = 2.61$, the standard error $s = 4.71$.

Is the gain significant?

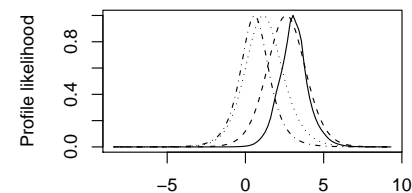
Intuitively: YES, if no two negative observations -8.4 and -6.0.

(a) Normal plot



QQ-plot: Quantile of $N(0, 1)$ vs Quantile of the empirical distribution.

(b) L_k likelihood, $k=1, 2, 4, 8$



The profile likelihood $l_k(\mu)$ is plotted against μ for $k = 1$ (solid), 2 (dashed), 4 (dotted), and 8 (dot-dashed).

(ii) Consider a generalised normal family

$$f_k(x|\mu, \sigma) = \frac{2^{-1-1/k}}{\Gamma(1+1/k)\sigma} \exp\left\{-\frac{1}{2}\left|\frac{x-\mu}{\sigma}\right|^k\right\},$$

which has the mean μ . When $k = 2$, it is $N(\mu, \sigma^2)$.

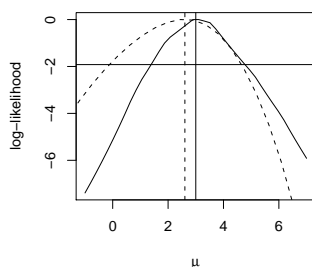
To find the profile likelihood of μ , the 'MLE' for σ is

$$\hat{\sigma}^k \equiv \hat{\sigma}(\mu)^k = \frac{k}{2n} \sum_{i=1}^n |X_i - \mu|^k.$$

Hence

$$l_k(\mu) = l_k(\mu, \hat{\sigma}) = -n \log \Gamma\left(1 + \frac{1}{k}\right) - n\left(1 + \frac{1}{k}\right) \log 2 - n \log \hat{\sigma} - \frac{n}{k}.$$

Fig.(b) shows the MLE $\hat{\mu} = \hat{\mu}(k)$ varies with respect to k . In fact $\hat{\mu}(k)$ increases as k decreases.



Parametric log-likelihood (solid curve) based on the DE distribution, and the empirical log-likelihood (dashed curve). (Both curves were shifted vertically by their own maximum values.)

If we use the distribution functions with $k = 1$ to fit the data, the p -value for the test is 0.03 – much more significant than that under the assumption of normal distribution.

(iii) The empirical likelihood ratio test statistic $2 \log T = 3.56$, which rejects H_0 with the p -value 0.04.

The 95% confidence interval is

$$\{\mu \mid \sum_{i=1}^{15} \log p_i(\mu) > -1.92 - 15 \log(15)\} = [0.17, 4.27].$$

The DE density is of the form $\frac{1}{2\sigma} e^{-|x-\mu|/\sigma}$. With μ fixed, the MLE for σ is $n^{-1} \sum_i |X_i - \mu|$. Hence the parametric log (profile) likelihood is

$$-n \log \sum_i |X_i - \mu|.$$

3. Empirical likelihood for random vectors

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d. random vectors from distribution F .

Similar to the univariate case, we assume

$$p_i = F(\mathbf{X}_i), \quad i = 1, \dots, n,$$

where $p_i \geq 0$ and $\sum_i p_i = 1$. The empirical likelihood is

$$L(p_1, \dots, p_n) = \prod_{i=1}^n p_i$$

Without any further constraints, the MELEs are

$$\hat{p}_i = 1/n, \quad i = 1, \dots, n.$$

3.1 EL for multivariate means

The profile empirical likelihood for $\boldsymbol{\mu} = E\mathbf{X}_1$ is

$$L(\boldsymbol{\mu}) = \max \left\{ \prod_{i=1}^n p_i \mid p_i \geq 0, \sum_{i=1}^n p_i = 1, \sum_{i=1}^n p_i \mathbf{X}_i = \boldsymbol{\mu} \right\} \equiv \prod_{i=1}^n p_i(\boldsymbol{\mu}),$$

where $p_i(\boldsymbol{\mu})$ is the MELE of p_i with the additional constraint $E\mathbf{X}_i = \boldsymbol{\mu}$. Define the ELR

$$T \equiv T(\boldsymbol{\mu}) = \frac{L(1/n, \dots, 1/n)}{L(\boldsymbol{\mu})} = 1 / \prod_{i=1}^n \{np_i(\boldsymbol{\mu})\}.$$

Theorem 3. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be $d \times 1$ i.i.d. with mean $\boldsymbol{\mu}$ and finite covariance matrix Σ and $|\Sigma| \neq 0$. Then as $n \rightarrow \infty$,

$$2 \log\{T(\boldsymbol{\mu})\} = -2 \sum_{i=1}^n \log\{np_i(\boldsymbol{\mu})\} \rightarrow \chi_d^2$$

in distribution.

(iv) *Bootstrap calibration.* Since (ii) and (iii) are based on an asymptotic result. When n is small and d is large, $\chi_{d,1-\alpha}^2$ may be replaced by *the $[B\alpha]$ -th largest value* among $2 \log T_1^*, \dots, 2 \log T_B^*$ which are computed as follows.

(a) Draw i.i.d. sample $\mathbf{X}_1^*, \dots, \mathbf{X}_n^*$ from the uniform distribution on $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$. Let

$$T^* = 1 / \prod_{i=1}^n \{np_i^*(\bar{\mathbf{X}})\},$$

where $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$, and $p_i^*(\boldsymbol{\mu})$ is obtained in the same manner as $p_i(\boldsymbol{\mu})$ with $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ replaced by $\{\mathbf{X}_1^*, \dots, \mathbf{X}_n^*\}$.

(b) Repeat (a) B times, denote the B values of T^* as T_1^*, \dots, T_B^* .

We may draw an \mathbf{X}^* from the uniform distribution on $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ as follows: draw $Z \sim U(0, 1)$, define $\mathbf{X}^* = \mathbf{X}_i$ if $Z \in [\frac{i-1}{n}, \frac{i}{n})$.

Remarks. (i) In the case that $|\Sigma| = 0$, there exists an integer $q < d$ for which, $\mathbf{X}_i = A\mathbf{Y}_i$ where \mathbf{Y}_i is a $q \times 1$ r.v. with $|\text{Var}(\mathbf{Y}_i)| \neq 0$, and A is a $d \times q$ constant matrix. The above theorem still holds with the limit distribution replaced by χ_q^2 .

(ii) The null hypothesis $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ will be rejected at the significance level α iff

$$\sum_{i=1}^n \log\{np_i(\boldsymbol{\mu}_0)\} \leq -0.5\chi_{d,1-\alpha}^2,$$

where $P\{\chi_d^2 \leq \chi_{d,1-\alpha}^2\} = 1 - \alpha$.

(iii) A $100(1 - \alpha)\%$ confidence region for $\boldsymbol{\mu}$ is

$$\left\{ \boldsymbol{\mu} \mid \sum_{i=1}^n \log\{np_i(\boldsymbol{\mu})\} \geq -0.5\chi_{d,1-\alpha}^2 \right\}.$$

Since the limiting distribution is free from the original distribution of $\{X_i\}$, we may draw X_i^* from any distribution $\{\pi_1, \dots, \pi_n\}$ instead of the uniform distribution used above. Of course now $p_i^*(\bar{\mathbf{X}})$ should be replaced by $p^*(\tilde{\boldsymbol{\mu}})$, where $\tilde{\boldsymbol{\mu}} = \sum_i \pi_i \mathbf{X}_i$.

(v) Computing $p_i(\boldsymbol{\mu})$.

Assumptions: $|\text{Var}(\mathbf{X}_i)| \neq 0$, and $\boldsymbol{\mu}$ is an inner point of the convex hull spanned by the observations, i.e.

$$\boldsymbol{\mu} \in \left\{ \sum_{i=1}^n p_i \mathbf{X}_i \mid p_i > 0, \sum_{i=1}^n p_i = 1 \right\}.$$

This ensures the solutions $p_i(\boldsymbol{\mu}) > 0$ exist.

We solve the problem in 3 steps:

1. Transform the constrained n -dim problem to a constrained d -dim problem.
2. Transform the constrained problem to an unconstrained problem.
3. Apply a Newton-Raphson algorithm.

Put

$$l(\mu) \equiv \log L(\mu) = \sum_{i=1}^n \log p_i(\mu)$$

$$= \max \left\{ \sum_{i=1}^n \log p_i \mid p_i \geq 0, \sum_{i=1}^n p_i = 1, \sum_{i=1}^n p_i \mathbf{X}_i = \mu \right\}.$$

Thus $M(\cdot)$ is a convex function on any connected sets satisfying

$$n - \lambda^\tau (\mathbf{X}_i - \mu) > 0, \quad i = 1, \dots, n. \quad (10)$$

Note. (10) and (8) together imply $\sum_i p_i(\mu) = 1$.

The original n -dimensional optimisation problem is equivalent to a d -dimensional problem of **minimising** $M(\lambda)$ subject to the constraints (10).

Let \mathcal{H}_λ be the set consisting all the values of λ satisfying

$$n - \lambda^\tau (\mathbf{X}_i - \mu) > 1, \quad i = 1, \dots, n.$$

Then \mathcal{H}_λ a convex set in R^d , which contains the minimiser of the convex function $M(\lambda)$. (See 'Note' above)

Unfortunately $M(\lambda)$ is not defined on the sets

$$\{\lambda \mid n - \lambda^\tau (\mathbf{X}_i - \mu) = 0\}, \quad i = 1, \dots, n.$$

Step 1:

Similar to Theorem 1, the LM method entails

$$p_i(\mu) = \frac{1}{n - \lambda^\tau (\mathbf{X}_i - \mu)}, \quad i = 1, \dots, n,$$

where λ is the solution of

$$\sum_{j=1}^n \frac{\mathbf{X}_j - \mu}{n - \lambda^\tau (\mathbf{X}_j - \mu)} = 0. \quad (8)$$

Hence

$$l(\mu) = - \sum_{i=1}^n \log \{n - \lambda^\tau (\mathbf{X}_i - \mu)\} \equiv M(\lambda). \quad (9)$$

Note $\frac{\partial}{\partial \lambda} M(\lambda) = 0$ leads to (8), and

$$\frac{\partial^2 M(\lambda)}{\partial \lambda \partial \lambda^\tau} = \sum_{i=1}^n \frac{(\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^\tau}{\{n - \lambda^\tau (\mathbf{X}_i - \mu)\}^2} > 0.$$

Step 2: We extend $M(\lambda)$ outside of the set \mathcal{H}_λ such that it is still a **convex function** on the whole R^d .

Define

$$\log_*(z) = \begin{cases} \log z & z \geq 1, \\ -1.5 + 2z - 0.5z^2 & z < 1. \end{cases}$$

It is easy to see that $\log_*(z)$ has two continuous derivatives on R .

Put $M_*(\lambda) = - \sum_{i=1}^n \log_* \{n - \lambda^\tau (\mathbf{X}_i - \mu)\}$. Then

- $M_*(\lambda) = M(\lambda)$ on \mathcal{H}_λ .
- $M_*(\lambda)$ is a convex function on whole R^d .

Hence, $M_*(\lambda)$ and $M(\lambda)$ share the same minimiser which is the solution of (8).

Step 3: We apply a Newton-Raphson algorithm to compute λ iteratively:

$$\lambda_{k+1} = \lambda_k - \{\dot{M}_*(\lambda_k)\}^{-1} \dot{M}_*(\lambda_k).$$

A convenient initial value would $\lambda_0 = 0$, corresponding to $p_i = 1/n$.

Remarks. (i) S-code "el.S", available from

www-stat.stanford.edu/~owen/empirical

calculates the empirical likelihood ratio

$$\sum_{i=1}^n \log\{np_i(\mu)\}$$

and other related quantities.

Theorem 4. Let X_1, \dots, X_n be $d \times 1$ i.i.d. r.v.s with mean μ_0 and $|\text{Var}(X_1)| \neq 0$. Let $\theta = h(\mu)$ be a smooth function from R^d to R^q ($q \leq d$), and $\theta_0 = h(\mu_0)$. We assume

$$|GG^T| \neq 0, \quad G = \frac{\partial \theta}{\partial \mu^T}.$$

For any $r > 0$, let

$$\mathcal{C}_{1,r} = \left\{ \mu \mid \sum_{i=1}^n \log\{np_i(\mu)\} \geq -0.5r \right\},$$

and

$$\mathcal{C}_{3,r} = \left\{ \theta_0 + G(\mu - \mu_0) \mid \mu \in \mathcal{C}_{1,r} \right\}.$$

Then as $n \rightarrow \infty$,

$$P\{\theta \in \mathcal{C}_{3,r}\} \rightarrow P(\chi_q^2 \leq r).$$

3.2 EL for smooth functions of means

Basic idea. Let Y_1, \dots, Y_n be i.i.d. random variables with variance σ^2 . Note

$$\sigma^2 = EY_i^2 - (EY_i)^2 = h(\mu),$$

where $\mu = EX_i$, and $X_i = (Y_i, Y_i^2)$. We may deduce a confidence interval for σ^2 from that of μ .

Remarks. (i) The idea of bootstrap calibration may be applied here too.

(ii) Under more conditions, $P\{\theta \in \mathcal{C}_{2,r}\} \rightarrow P(\chi_q^2 \leq r)$, where

$$\mathcal{C}_{2,r} = \left\{ h(\mu) \mid \mu \in \mathcal{C}_{1,r} \right\}.$$

$\mathcal{C}_{2,r}$ is a practical feasible confidence set, while $\mathcal{C}_{3,r}$ is not since μ_0 and θ_0 are unknown in practice. Note for μ close to μ_0 ,

$$\theta_0 + G(\mu - \mu_0) \approx h(\mu).$$

(iii) In general, $P\{\mu \in \mathcal{C}_{1,r}\} \leq P\{\theta \in \mathcal{C}_{2,r}\}$.

(By Theorem 3, $P\{\mu \in \mathcal{C}_{1,r}\} \rightarrow P(\chi_d^2 \leq r)$)

(iv) The **profile empirical likelihood function** of θ is

$$L(\theta) = \max \left\{ \prod_{i=1}^n p_i(\mu) \mid h(\mu) = \theta \right\}$$

$$= \max \left\{ \prod_{i=1}^n p_i \mid h\left(\sum_{i=1}^n p_i X_i\right) = \theta, p_i \geq 0, \sum_{i=1}^n p_i = 1 \right\},$$

which may be calculated directly using the Lagrange multiplier method. The computation is more involved for nonlinear $h(\cdot)$.

Example 4. S&P500 stock index in 17.8.1999 — 17.8.2000 (256 trading days)

Let Y_i be the price on the i -th day,

$$X_i = \log(Y_i/Y_{i-1}) \approx (Y_i - Y_{i-1})/Y_{i-1},$$

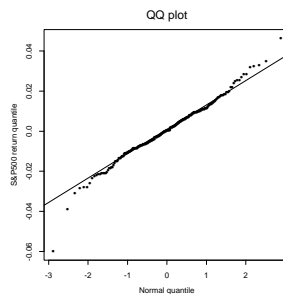
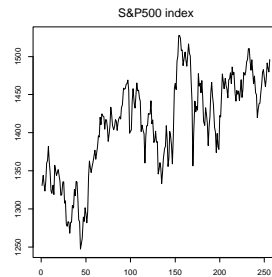
which is the return, i.e. the percentage of the change on the i -th day.

By treating X_i i.i.d., we construct confidence intervals for the annual volatility

$$\sigma = \{255 \text{Var}(X_i)\}^{1/2}.$$

The simple point-estimator is

$$\hat{\sigma} = \left\{ \frac{255}{255} \sum_{i=1}^{255} (X_i - \bar{X})^2 \right\}^{1/2} = 0.2116.$$



The 95% confidence intervals are:

Method	C.I.
EL	[0.1895, 0.2422]
Normal	[0.1950, 0.2322]

The EL confidence interval is **41.67%** wider than the interval based on normal distribution, which reflects the fact that **the returns have heavier tails**.

4. Estimating equations

4.1 Estimation via estimating equations

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d. from a distribution F . We are interested in some characteristic $\boldsymbol{\theta} \equiv \boldsymbol{\theta}(F)$, which is determined by equation

$$E\{m(\mathbf{X}_1, \boldsymbol{\theta})\} = 0,$$

where $\boldsymbol{\theta}$ is $q \times 1$ vector, m is a $s \times 1$ vector-valued function.

For example,

$$\theta = EX_1 \text{ if } m(x, \theta) = x - \theta,$$

$$\theta = E(X_1^k) \text{ if } m(x, \theta) = x^k - \theta,$$

$$\theta = P(X_1 \in A) \text{ if } m(x, \theta) = I(x \in A) - \theta,$$

$$\theta \text{ is the } \alpha\text{-quantile if } m(x, \theta) = I(x \leq \theta) - \alpha.$$

Example 5. Let $\{(X_i, Y_i), i = 1, \dots, n\}$ be a random sample. Find a set of estimating equations for estimating $\gamma \equiv \text{Var}(X_1)/\text{Var}(Y_1)$.

In order to estimate γ , we need to estimate $\mu_x = E(X_1)$, $\mu_y = E(Y_1)$ and $\sigma_y^2 = \text{Var}(Y_1)$. Put $\boldsymbol{\theta}^T = (\mu_x, \mu_y, \sigma_y^2, \gamma)$, and

$$m_1(X, Y, \boldsymbol{\theta}) = X - \mu_x, \quad m_2(X, Y, \boldsymbol{\theta}) = Y - \mu_y,$$

$$m_3(X, Y, \boldsymbol{\theta}) = (Y - \mu_y)^2 - \sigma_y^2,$$

$$m_4(X, Y, \boldsymbol{\theta}) = (X - \mu_x)^2 - \sigma_y^2 \gamma,$$

and $\mathbf{m} = (m_1, m_2, m_3, m_4)^T$. Then $E\{\mathbf{m}(X_i, Y_i, \boldsymbol{\theta})\} = 0$, leading to the estimating equation

$$\frac{1}{n} \sum_{i=1}^n \mathbf{m}(X_i, Y_i, \boldsymbol{\theta}) = 0,$$

the solution of the above equation is an estimator $\hat{\boldsymbol{\theta}}$ for $\boldsymbol{\theta}$.

Remark. Estimating equation method does not facilitate hypothesis tests and interval estimation for $\boldsymbol{\theta}$.

A natural estimator for $\boldsymbol{\theta}$ is determined by *the estimating equation*

$$\frac{1}{n} \sum_{i=1}^n m(\mathbf{X}_i, \hat{\boldsymbol{\theta}}) = 0. \quad (11)$$

Obviously, in case F is in a parametric family and m is the score function, $\hat{\boldsymbol{\theta}}$ is the ordinary MLE.

Determined case $q = s$: $\hat{\boldsymbol{\theta}}$ may be uniquely determined by (11)

Underdetermined case $q > s$: the solutions of (11) may form a $(q - s)$ -dimensional set

Overdetermined case $q < s$: (11) may not have an exact solution, approximating solutions are sought. One such an example is so-called *the generalised method of moments estimation* which is very popular in Econometrics.

4.2 EL for estimating equations

Aim: construct statistical tests and confidence intervals for $\boldsymbol{\theta}$

The profile empirical likelihood function of $\boldsymbol{\theta}$:

$$L(\boldsymbol{\theta}) = \max \left\{ \prod_{i=1}^n p_i \mid \sum_{i=1}^n p_i m(\mathbf{X}_i, \boldsymbol{\theta}) = 0, p_i \geq 0, \sum_{i=1}^n p_i = 1 \right\}$$

The following theorem follows from Theorem 2 immediately.

Theorem 5. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d., $m(\mathbf{x}, \boldsymbol{\theta})$ be an $s \times 1$ vector-valued function. Suppose

$$E\{m(\mathbf{X}_1, \boldsymbol{\theta}_0)\} = 0, \quad \left| \text{Var}\{m(\mathbf{X}_1, \boldsymbol{\theta}_0)\} \right| \neq 0.$$

Then as $n \rightarrow \infty$,

$$-2 \log\{L(\boldsymbol{\theta}_0)\} - 2n \log n \rightarrow \chi_s^2$$

in distribution.

The theorem above applies in all *determined, underdetermined and overdetermined cases*.

Remarks (i) In general $L(\theta)$ can be calculated using the method for EL for multivariate means in §3.1, treating $m(\mathbf{X}_i, \theta)$ as a random vector.

(ii) For $\theta = \hat{\theta}$ which is the solution of

$$\frac{1}{n} \sum_{i=1}^n m(\mathbf{X}_i, \hat{\theta}) = 0,$$

$$L(\hat{\theta}) = (1/n)^n.$$

(iii) For θ determined by $E\{m(\mathbf{X}_1, \theta)\} = 0$, we will reject the null hypothesis $H_0 : \theta = \theta_0$ iff

$$\log\{L(\theta_0)\} + n \log n \leq -0.5\chi_{s,1-\alpha}^2.$$

(iii) An $(1-\alpha)$ confidence set for θ determined by $E\{m(\mathbf{X}_1, \theta)\} = 0$ is

$$\{\theta \mid \log\{L(\theta)\} + n \log n > -0.5\chi_{s,1-\alpha}^2\}$$

Let

$$L(\theta_\alpha) = \max \left\{ \prod_{i=1}^n p_i \mid \sum_{i=1}^n p_i I(X_i \leq \theta_\alpha) = \alpha, \right. \\ \left. p_i \geq 0, \sum_{i=1}^n p_i = 1 \right\}.$$

An $(1-\beta)$ confidence interval for the α quantile is

$$\Theta_\alpha = \{\theta_\alpha \mid \log\{L(\theta_\alpha)\} > -n \log n - 0.5\chi_{1,1-\beta}^2\}.$$

Note $L(\hat{\theta}_\alpha) = (1/n)^n \geq L(\theta_\alpha)$ for any θ_α . It is always true that $\hat{\theta}_\alpha \in \Theta_\alpha$.

Example 6. (Confidence intervals for quantiles)

Let X_1, \dots, X_n be i.i.d. For a given $\alpha \in (0, 1)$, let

$$m(x, \theta_\alpha) = I(x \leq \theta_\alpha) - \alpha.$$

Then $E\{m(X_i, \theta_\alpha)\} = 0$ implies θ_α is the α quantile of the distribution of X_i . We assume the true value of θ_α is between $X_{(1)}$ and $X_{(n)}$.

The estimating equation

$$\sum_{i=1}^n m(X_i, \hat{\theta}_\alpha) = \sum_{i=1}^n I(X_i \leq \hat{\theta}_\alpha) - n\alpha = 0$$

entails

$$\hat{\theta}_\alpha = X_{(n\alpha)},$$

where $X_{(i)}$ denotes the i -th smallest value among X_1, \dots, X_n . We assume $n\alpha$ is an integer to avoid insignificant (for large n , e.g. $n = 100$) technical details.

In fact $L(\theta_\alpha)$ can be computed explicitly as follows.

Let $r = r(\theta_\alpha)$ be the integer for which

$$X_{(i)} \leq \theta_\alpha \text{ for } i = 1, \dots, r, \text{ and}$$

$$X_{(i)} > \theta_\alpha \text{ for } i = r + 1, \dots, n.$$

Thus

$$L(\theta_\alpha) = \max \left\{ \prod_{i=1}^n p_i \mid p_i \geq 0, \sum_{i=1}^r p_i = \alpha, \sum_{i=r+1}^n p_i = 1 - \alpha \right\} \\ = (\alpha/r)^r \{(1-\alpha)/(n-r)\}^{n-r}.$$

Hence

$$\begin{aligned}\Theta_\alpha &= \{\theta_\alpha \mid \log\{L(\theta_\alpha)\} > -n \log n - 0.5\chi_{1,1-\alpha}^2\} \\ &= \left\{ \theta_\alpha \mid r \log \frac{n\alpha}{r} + (n-r) \log \frac{n(1-\alpha)}{n-r} > -0.5\chi_{1,1-\alpha}^2 \right\},\end{aligned}$$

which can also be derived directly based on a likelihood ratio test for a binomial distribution.

5.1 From global fitting to local fitting

Consider linear regression model

$$Y = X_1\beta_1 + \dots + X_d\beta_d + \varepsilon, \quad (12)$$

where $\varepsilon \sim (0, \sigma^2)$.

This model is *linear wrt unknown coefficients* β_1, \dots, β_d as the variable X_1, \dots, X_d may be

- quantitative inputs
- transformations of quantitative inputs, such as log, square-root etc
- interactions between variables, e.g. $X_3 = X_1X_2$
- basis expansions, such as $X_2 = X_1^2, X_3 = X_1^3, \dots$
- numeric or “dummy” coding of the levels if qualitative inputs

5. Empirical likelihood for estimating conditional distributions

References on kernel regression:

- Simonoff, J. S. (1996). *Smoothing Methods in Statistics*. Springer, New York.
- Wand, M.P. and Jones, M.C. (1995). *Kernel Smoothing*. Chapman and Hall, London.

References on nonparametric estimation for distribution functions:

- Hall, P., Wolff, R.C.L. and Yao, Q. (1999). Methods for estimating a conditional distribution function. *Journal of the American Statistical Association*, **94**, 154-163.
- Fan, J. and Yao, Q. (2003). *Nonlinear Time Series: Nonparametric and Parametric Methods*. Springer, New York. Sections 10.3 (also Section 6.5).

Put $\beta = (\beta_1, \dots, \beta_d)^\tau$

With observations $\{(Y_i, \mathbf{X}_i), 1 \leq i \leq n\}$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{id})^\tau$, the LSE minimises

$$\sum_{i=1}^n (Y_i - \mathbf{X}_i^\tau \beta)^2, \quad (13)$$

resulting to

$$\hat{\beta} = (\mathbf{X}^\tau \mathbf{X})^{-1} \mathbf{X}^\tau \mathbf{Y},$$

where $\mathbf{Y} = (Y_1, \dots, Y_n)^\tau$, and $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^\tau$ is an $n \times d$ matrix.

The fitted model is

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\beta}.$$

This is a **global** fitting, since the model is assumed to be true everywhere in the sample space and the estimator $\hat{\beta}$ is obtained using all the available data.