

Estimation of deformation between distributions by minimal Wasserstein distance

Hélène Lescornel

Institut de Mathématiques de Toulouse

Colloque Jeunes Probabilistes et Statisticiens - 7/04/2014

Summary

Introduction

The model

- Statistical framework

- The estimators

Consistency

- M-estimation

- Result

Convergence in distribution

- New framework

- Idea of proof

Examples

Summary

Introduction

The model

Consistency

Convergence in distribution

Examples

- We observe a random variable and a deformation of this variable

$$\begin{cases} \varepsilon \\ X = \varphi(\varepsilon) \end{cases}$$

⇒ Random experiments with some variability : φ .



- We observe a random variable and a deformation of this variable

$$\begin{cases} \varepsilon \\ X = \varphi(\varepsilon) \end{cases}$$

⇒ Random experiments with some variability : φ .

- How to extract information? : estimation of the deformation, study of a mean distribution if several deformations are observed...

Warped curves

- Dynamic Time Warping. Sakoe-Chiba-[1978]

Align two signals $(f(i))_{1 \leq i \leq N}$ and $(g(j))_{1 \leq j \leq M}$ by a time axis re normalization.

Idea : to consider some "warping operators" between $1 \leq i \leq N$ and $1 \leq j \leq M$.

↪ Minimize a cost

$$C(w, f, g) = \sum_{(i,j) \in w} (f(i) - g(j))^2.$$

Warped curves

- Dynamic Time Warping. Sakoe-Chiba-[1978]

Align two signals $(f(i))_{1 \leq i \leq N}$ and $(g(j))_{1 \leq j \leq M}$ by a time axis re normalization.

Idea : to consider some "warping operators" between $1 \leq i \leq N$ and $1 \leq j \leq M$.

↪ Minimize a cost

$$C(w, f, g) = \sum_{(i,j) \in w} (f(i) - g(j))^2.$$

- Extension to warped curves in a regression scheme : **Wang-Gasser-[1999], Gamboa-Loubes-Maza-[2007]**. Different cost functions.

Deformation of distributions

- Observations $X_j = \varphi_j(\varepsilon)$, $1 \leq j \leq J$ where φ_j are realizations of a random process.

Estimation of a mean distribution using the quantile functions F_j^{-1} in **Gallòn-Loubes-Maza-[2013]**.

Deformation of distributions

- Observations $X_j = \varphi_j(\varepsilon)$, $1 \leq j \leq J$ where φ_j are realizations of a random process.
Estimation of a mean distribution using the quantile functions F_j^{-1} in **Gallòn-Loubes-Maza-[2013]**.
- Test for a parametric relationship between two distributions in terms of quantile functions $F^{-1} = \mathcal{F}(G^{-1}, \theta)$. Test statistic based on the L^2 norm between quantile functions in **Freitag-Munk-[2005]**.

Model studied

Semi parametric framework

$$\begin{cases} \varepsilon \\ X = \varphi_{\theta^*}(\varepsilon) \end{cases}$$

- shape of the deformation φ assumed **known**,
- deformation parameter θ^* and template measure μ of ε to estimate.

Model studied

Semi parametric framework

$$\begin{cases} \varepsilon \\ X = \varphi_{\theta^*}(\varepsilon) \end{cases}$$

- shape of the deformation φ assumed **known**,
- deformation parameter θ^* and template measure μ of ε to estimate.

Idea: **Align** the distribution of X on the distribution of ε .

\rightsquigarrow Study of $Z(\theta) = \varphi_{\theta}^{-1}(X) = \varphi_{\theta}^{-1} \circ \varphi_{\theta^*}(\varepsilon)$.



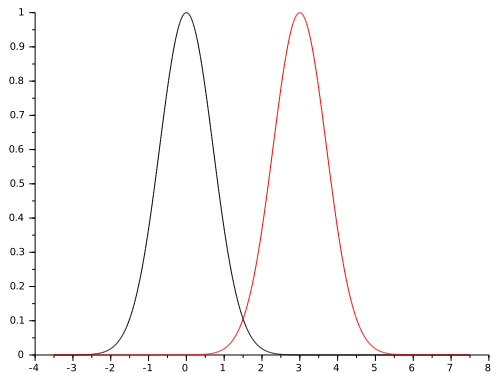
Parallel with warped curves

$$\varphi_{\theta}(t) = t + \theta$$



Parallel with warped curves

$$\varphi_{\theta}(t) = t + \theta$$



We have $Z(\theta^*) = \varepsilon$.

→ Align the distribution of $Z(\theta)$ on the distribution of ε by varying θ . Minimization of a D.T.W. criterion.



We have $Z(\theta^*) = \varepsilon$.

→ Align the distribution of $Z(\theta)$ on the distribution of ε by varying θ . Minimization of a D.T.W. criterion.

→ Which cost function? Requires a distance between probabilities. Utilization of the **Wasserstein distance** related to problems of mass transport.

Summary

Introduction

The model

Statistical framework

The estimators

Consistency

Convergence in distribution

Examples



Presentation of the model

Observations :

$$\begin{cases} \varepsilon_{i1} & 1 \leq i \leq n \\ X_i = \varphi_{\theta^*}(\varepsilon_{i2}) & 1 \leq i \leq n \end{cases} \quad (1)$$

- structure : $(\varepsilon_{ij})_{\substack{1 \leq i \leq n \\ j=1,2}}$ i.i.d. following the law μ ,
- deformation parameter : θ^* where $\theta^* \in \Theta \subset \mathbb{R}^d$,
- deformation function : $\varphi_\theta :]a; b[\rightarrow]c; d[$ invertible for $\theta \in \Theta$,
- Empirical distribution of the i.i.d. sample $(\varepsilon_{i1})_{1 \leq i \leq n} : \mu_1^n$.



For $\theta \in \Theta$ we define

$$Z_i(\theta) = \varphi_\theta^{-1}(X_i) = \varphi_\theta^{-1} \circ \varphi_{\theta^*}(\varepsilon_{i2})$$

- Distribution of $Z_1(\theta)$: $\mu_\star(\theta) = \mu \circ \varphi_{\theta^*}^{-1} \circ \varphi_\theta$,
- Empirical distribution of the i.i.d. sample $(Z_1(\theta), \dots, Z_n(\theta))$:
 $\mu_\star^n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Z_i(\theta)}$.



For $\theta \in \Theta$ we define

$$Z_i(\theta) = \varphi_\theta^{-1}(X_i) = \varphi_\theta^{-1} \circ \varphi_{\theta^*}(\varepsilon_{i2})$$

- Distribution of $Z_1(\theta)$: $\mu_\star(\theta) = \mu \circ \varphi_{\theta^*}^{-1} \circ \varphi_\theta$,
- Empirical distribution of the i.i.d. sample $(Z_1(\theta), \dots, Z_n(\theta))$:
 $\mu_\star^n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{Z_i(\theta)}$.

\Rightarrow Recover θ^\star by minimizing the energy needed to align the distribution $\mu_\star(\theta)$ on μ : $\mu_\star(\theta^\star) = \mu$.

Wasserstein distance to quantify the alignment.



Wasserstein distance

Set $\mathcal{W}_2(\mathbb{R}) = \{P \text{ probability on } \mathbb{R}, \int_{\mathbb{R}} x^2 dP < \infty\}$.

For $P, Q \in \mathcal{W}_2(\mathbb{R})$ with respective distribution function F and G , their Wasserstein distance is

$$W_2^2(P, Q) = \int_0^1 (F^{-1}(t) - G^{-1}(t))^2 dt. \quad (2)$$



Wasserstein distance

Set $\mathcal{W}_2(\mathbb{R}) = \{P \text{ probability on } \mathbb{R}, \int_{\mathbb{R}} x^2 dP < \infty\}$.

For $P, Q \in \mathcal{W}_2(\mathbb{R})$ with respective distribution function F and G , their Wasserstein distance is

$$W_2^2(P, Q) = \int_0^1 (F^{-1}(t) - G^{-1}(t))^2 dt. \quad (2)$$

If P and Q are defined on more general metric space (S, d) with a moment of order 2 :

$$W_2^2(P, Q) = \inf_{X \sim P, Y \sim Q} \mathbb{E} [d(X, Y)^2]$$

Wasserstein distance

Set $\mathcal{W}_2(\mathbb{R}) = \{P \text{ probability on } \mathbb{R}, \int_{\mathbb{R}} x^2 dP < \infty\}$.

For $P, Q \in \mathcal{W}_2(\mathbb{R})$ with respective distribution function F and G , their Wasserstein distance is

$$W_2^2(P, Q) = \int_0^1 (F^{-1}(t) - G^{-1}(t))^2 dt. \quad (2)$$

If P and Q are defined on more general metric space (S, d) with a moment of order 2 :

$$W_2^2(P, Q) = \inf_{X \sim P, Y \sim Q} \mathbb{E} [d(X, Y)^2]$$

- Set $X_n \sim P_n, X \sim P$.

$$\text{Then } W_2(P_n, P) \rightarrow 0 \iff \begin{cases} X_n \rightarrow X \\ \mathbb{E} [X_n^2] \rightarrow \mathbb{E} [X^2] \end{cases}$$



The estimators

To quantify the alignment of the measures, consider the criterion

$$M(\theta) = W_2^2(\mu_\star(\theta), \mu) \quad (3)$$



The estimators

To quantify the alignment of the measures, consider the criterion

$$M(\theta) = W_2^2(\mu_\star(\theta), \mu) \quad (3)$$

Characterisation of θ^\star : $\min_{\Theta} M = 0 = M(\theta^\star)$.



The estimators

To quantify the alignment of the measures, consider the criterion

$$M(\theta) = W_2^2(\mu_\star(\theta), \mu) \quad (3)$$

Characterisation of θ^\star : $\min_{\Theta} M = 0 = M(\theta^\star)$.

→ Empirical version :

$$M_n(\theta) = W_2^2(\mu_\star^n(\theta), \mu_1^n)$$



with the order statistics

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n [Z_{(i)}(\theta) - \varepsilon_{(i)1}]^2. \quad (4)$$



with the order statistics

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n [Z_{(i)}(\theta) - \varepsilon_{(i)1}]^2. \quad (4)$$

Leads to the M-estimator for the deformation parameters

Estimator of θ^*

$$\hat{\theta}^n \in \operatorname{argmin}_{\theta \in \Theta} M_n(\theta). \quad (5)$$



with the order statistics

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n [Z_{(i)}(\theta) - \varepsilon_{(i)1}]^2. \quad (4)$$

Leads to the M-estimator for the deformation parameters

Estimator of θ^*

$$\hat{\theta}^n \in \operatorname{argmin}_{\theta \in \Theta} M_n(\theta). \quad (5)$$

We have $\varphi_{\hat{\theta}^n}^{-1}(X_i) = \varphi_{\hat{\theta}^n}^{-1} \circ \varphi_{\theta^*}(\varepsilon_{i2}) \approx \varepsilon_{i2}$, following the unknown law μ .



with the order statistics

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n [Z_{(i)}(\theta) - \varepsilon_{(i)1}]^2. \quad (4)$$

Leads to the M-estimator for the deformation parameters

Estimator of θ^*

$$\hat{\theta}^n \in \operatorname{argmin}_{\theta \in \Theta} M_n(\theta). \quad (5)$$

We have $\varphi_{\hat{\theta}^n}^{-1}(X_i) = \varphi_{\hat{\theta}^n}^{-1} \circ \varphi_{\theta^*}(\varepsilon_{i2}) \approx \varepsilon_{i2}$, following the unknown law μ .

- Plugg-in estimator of μ

$$\hat{\mu}^n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\varphi_{\hat{\theta}^n}^{-1}(X_i)} \quad (6)$$

Idea : "increase" the size of $(\varepsilon_{i1})_{1 \leq i \leq n} \rightsquigarrow \tilde{\mu}^n = \frac{1}{2}(\hat{\mu}^n + \mu_1^n)$

Summary

Introduction

The model

Consistency

M-estimation

Result

Convergence in distribution

Examples



Principle of M-estimation

→ Convergence criterion for estimators defined as minimizers of a random functional M_n on a set Θ .

$$\begin{array}{ccc}
 M_n(\theta) & \xrightarrow{n \rightarrow \infty} & M(\theta) \\
 \min \uparrow & & \min \uparrow \\
 \hat{\theta}^n & \xrightarrow{?} & \theta^*
 \end{array}$$

Principle of M-estimation

→ Convergence criterion for estimators defined as minimizers of a random functional M_n on a set Θ .

$$\begin{array}{ccc}
 M_n(\theta) & \xrightarrow{n \rightarrow \infty} & M(\theta) \\
 \min \uparrow & & \min \uparrow \\
 \hat{\theta}^n & \xrightarrow{?} & \theta^*
 \end{array}$$

Consistency criterion

- If M is a deterministic function,
- $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{n \rightarrow \infty} 0$ in probability,
- $\forall \delta > 0 \quad \inf_{\theta \in \Theta \cap B(\theta^*, \delta)^c} M(\theta) > M(\theta^*)$

then $\hat{\theta}^n \xrightarrow{n \rightarrow \infty} \theta^*$ in probability.



Assumptions

- Laws considered are defined on subsets of \mathbb{R} and $\forall \theta \in \Theta$, $\mu_{\star}(\theta) \in \mathcal{W}_2(\mathbb{R})$.
 \Rightarrow Computation of the Wasserstein distance.



Assumptions

- Laws considered are defined on subsets of \mathbb{R} and $\forall \theta \in \Theta$, $\mu_{\star}(\theta) \in \mathcal{W}_2(\mathbb{R})$.
 \Rightarrow Computation of the Wasserstein distance.
- Regularity C^1 of $\varphi_{\theta}^{-1}(x)$ with respect to $\theta \in \Theta$.

Assumptions

- Laws considered are defined on subsets of \mathbb{R} and $\forall \theta \in \Theta$, $\mu_\star(\theta) \in \mathcal{W}_2(\mathbb{R})$.
 \Rightarrow Computation of the Wasserstein distance.
- Regularity C^1 of $\varphi_\theta^{-1}(x)$ with respect to $\theta \in \Theta$. The family $\{\partial\varphi_\theta^{-1}(\cdot)\}_{\theta \in \Theta}$ has an envelop in $L^2(X)$:

$$\sup_{\theta \in \Theta} \|\partial\varphi_\theta^{-1}(x)\| \leq H(x) \text{ with } H \in L^2(X)$$

\Rightarrow Control the distance between $\mu_\star(\theta^1)$ and $\mu_\star(\theta^2)$ for $\theta^1, \theta^2 \in \Theta$.

Assumptions

- Laws considered are defined on subsets of \mathbb{R} and $\forall \theta \in \Theta$, $\mu_\star(\theta) \in \mathcal{W}_2(\mathbb{R})$.
 \Rightarrow Computation of the Wasserstein distance.
- Regularity C^1 of $\varphi_\theta^{-1}(x)$ with respect to $\theta \in \Theta$. The family $\{\partial\varphi_\theta^{-1}(\cdot)\}_{\theta \in \Theta}$ has an envelop in $L^2(X)$:

$$\sup_{\theta \in \Theta} \|\partial\varphi_\theta^{-1}(x)\| \leq H(x) \text{ with } H \in L^2(X)$$

\Rightarrow Control the distance between $\mu_\star(\theta^1)$ and $\mu_\star(\theta^2)$ for $\theta^1, \theta^2 \in \Theta$.

- Θ compact and convex subset of \mathbb{R}^d .
 \Rightarrow Uniform convergence and Taylor expansion.

Assumptions

- Laws considered are defined on subsets of \mathbb{R} and $\forall \theta \in \Theta$, $\mu_\star(\theta) \in \mathcal{W}_2(\mathbb{R})$.
 \Rightarrow Computation of the Wasserstein distance.
- Regularity C^1 of $\varphi_\theta^{-1}(x)$ with respect to $\theta \in \Theta$. The family $\{\partial\varphi_\theta^{-1}(\cdot)\}_{\theta \in \Theta}$ has an envelop in $L^2(X)$:

$$\sup_{\theta \in \Theta} \|\partial\varphi_\theta^{-1}(x)\| \leq H(x) \text{ with } H \in L^2(X)$$

\Rightarrow Control the distance between $\mu_\star(\theta^1)$ and $\mu_\star(\theta^2)$ for $\theta^1, \theta^2 \in \Theta$.

- Θ compact and convex subset of \mathbb{R}^d .
 \Rightarrow Uniform convergence and Taylor expansion.
- Identifiability condition : for all $\theta \neq \theta^\star$, $\varphi_\theta^{-1} \circ \varphi_{\theta^\star} \neq Id$ on a set of positive μ -measure.
 \Rightarrow Uniqueness of the minimizer of the function M .

Consistency results

Deformation estimator

$$\hat{\theta}^n \in \operatorname{argmin}_{\theta \in \Theta} M_n(\theta) :$$

Theorem

Under previous assumptions $\hat{\theta}^n$ converges in probability to θ^ .*

Measure estimator

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\varphi_{\hat{\theta}^n}^{-1}}(X_i)$$

Theorem

Under previous assumptions

$$W_2(\hat{\mu}_n, \mu) \xrightarrow{n \rightarrow \infty} 0 \text{ in probability.}$$

Summary

Introduction

The model

Consistency

Convergence in distribution

New framework

Idea of proof

Examples

Assumptions

In addition to the previous assumptions, we assume

- More regularity : φ^{-1} is C^2 with respect to its two variables (θ, x) .



Assumptions

In addition to the previous assumptions, we assume

- More regularity : φ^{-1} is C^2 with respect to its two variables (θ, x) .
- The distribution of X has a compact support with distribution function $F_\star \in C^1$. We assume $F'_\star := f_\star > 0$ on its support.

\Rightarrow The distribution function F associated with the law μ (law of ε) has a compact support and is C^1 with $F' = f > 0$.

Result

Set $\Phi = \int_0^1 \partial \varphi_{\theta^*}^{-1} (F_{\star}^{-1}(t))^2 dt \in \mathbb{R}^{d \times d}$.

Theorem

Under previous assumptions and if Φ is invertible, then

$$\sqrt{n} (\hat{\theta}^n - \theta^*) \rightharpoonup \Phi^{-1} \int_0^1 \frac{\partial \varphi_{\theta^*}^{-1} (F_{\star}^{-1}(t))}{f(F_{\star}^{-1}(t))} [\mathbb{G}_2(t) - \mathbb{G}_1(t)] dt$$

where \mathbb{G}_1 and \mathbb{G}_2 are independent standard Brownian bridges.



Idea of proof

→ Remains to study $\Psi(F^n, F_\star^n)$ where F^n (resp. F_\star^n) is the empirical distribution function associated with the sample $(\varepsilon_{i1})_{1 \leq i \leq n}$ (resp. $(X_i)_{1 \leq i \leq n}$).

Idea of proof

→ Remains to study $\Psi(F^n, F_\star^n)$ where F^n (resp. F_\star^n) is the empirical distribution function associated with the sample $(\varepsilon_{i1})_{1 \leq i \leq n}$ (resp. $(X_i)_{1 \leq i \leq n}$).

Convergence of the empirical distribution functions :

Theorem (Donsker)

If Y_1, \dots, Y_n are i.i.d. random variables with distribution function F and empirical distribution function F_n , the sequence $\sqrt{n}(F_n - F)$ converges in law in \mathbb{S} , the space of function cadlag on $\bar{\mathbb{R}}$ endowed with the norm $\|\cdot\|_\infty$ to $\mathbb{G} \circ F$ where \mathbb{G} is a standard Brownian bridge.

Idea of proof

→ Remains to study $\Psi(F^n, F_\star^n)$ where F^n (resp. F_\star^n) is the empirical distribution function associated with the sample $(\varepsilon_{i1})_{1 \leq i \leq n}$ (resp. $(X_i)_{1 \leq i \leq n}$).

Convergence of the empirical distribution functions :

Theorem (Donsker)

If Y_1, \dots, Y_n are i.i.d. random variables with distribution function F and empirical distribution function F_n , the sequence $\sqrt{n}(F_n - F)$ converges in law in \mathbb{S} , the space of function cadlag on $\bar{\mathbb{R}}$ endowed with the norm $\|\cdot\|_\infty$ to $\mathbb{G} \circ F$ where \mathbb{G} is a standard Brownian bridge.

⇒ Application of a **Delta-method**.

Summary

Introduction

The model

Consistency

Convergence in distribution

Examples

Examples

- Example 1 : Translation model

$$\varphi_{\theta}(x) = x + \theta$$

$\Rightarrow \mu \in \mathcal{W}_2(\mathbb{R})$, and $\Theta \subset \mathbb{R}$ compact interval.

Examples

- Example 1 : Translation model

$$\varphi_{\theta}(x) = x + \theta$$

$\Rightarrow \mu \in \mathcal{W}_2(\mathbb{R})$, and $\Theta \subset \mathbb{R}$ compact interval.

$$\begin{cases} \varepsilon_{i1} \\ X_i = \theta^* + \varepsilon_{i2} \end{cases} \quad 1 \leq i \leq n$$

$$\hat{\theta}^n = \frac{1}{n} \sum_{i=1}^n X_i - \frac{1}{n} \sum_{i=1}^n \varepsilon_{i1} = \theta^* - \left[\frac{1}{n} \sum_{i=1}^n \varepsilon_{i2} - \frac{1}{n} \sum_{i=1}^n \varepsilon_{i1} \right]$$

- Example 2 : Logit model

$$\varphi_{\theta}(x) = \frac{1}{1 + \exp(\theta x)}$$

$\Rightarrow \mu \in \mathcal{W}_2(\mathbb{R})$ and Θ compact interval of $] - \infty; 0[$.

- Example 2 : Logit model

$$\varphi_{\theta}(x) = \frac{1}{1 + \exp(\theta x)}$$

$\Rightarrow \mu \in \mathcal{W}_2(\mathbb{R})$ and Θ compact interval of $] - \infty; 0[$.

$$\begin{cases} \varepsilon_{i1} \\ X_i = \frac{1}{1 + \theta^* \varepsilon_{i2}} \end{cases} \quad 1 \leq i \leq n$$

$$\hat{\theta}^n = \frac{\sum_{i=1}^n \ln \left(\frac{1 - X_{(i)}}{X_{(i)}} \right)^2}{\sum_{i=1}^n \ln \left(\frac{1 - X_{(i)}}{X_{(i)}} \right) \varepsilon_{(i)1}} = \theta^* \frac{\sum_{i=1}^n \varepsilon_{(i)2}^2}{\sum_{i=1}^n \varepsilon_{(i)2} \varepsilon_{(i)1}}$$



- Example 3: Location/scale model

$$\varphi_{\theta}(x) = \theta_2 x + \theta_1$$

$\Rightarrow \mu \in \mathcal{W}_2(\mathbb{R})$ and Θ compact in $\mathbb{R} \times]0; +\infty[$.

- Example 3: Location/scale model

$$\varphi_{\theta}(x) = \theta_2 x + \theta_1$$

$\Rightarrow \mu \in \mathcal{W}_2(\mathbb{R})$ and Θ compact in $\mathbb{R} \times]0; +\infty[$.

→ Scale model

$$\begin{cases} \varepsilon_{i1} \\ X_i = \theta^* \varepsilon_{i2} \end{cases} \quad 1 \leq i \leq n$$

$$\hat{\theta}^n = \frac{\sum_{i=1}^n X_{(i)}^2}{\sum_{i=1}^n X_{(i)} \varepsilon_{(i)1}} = \theta^* \frac{\sum_{i=1}^n \varepsilon_{(i)2}^2}{\sum_{i=1}^n \varepsilon_{(i)2} \varepsilon_{(i)1}}$$

Bibliography

Gamboa-Loubes-Maza-[2007] : F. Gamboa, J.-M. Loubes, and E. Maza. Semi-parametric estimation of shifts. Electron. J. Stat., 1 :616-640, 2007.

Gallòn-Loubes-Maza-[2011] : S. Gallòn, J-M Loubes, E. Maza, Statistical Properties of the Quantile Normalization Method for Density Curve Alignment. Technical report, May 2011.

Vimond-[2010] : M. Vimond, Efficient estimation for a subclass of shape invariant models. Ann. Statist., 38(3):1885-1912, 2010.

Wang-Gasser-[1999] : K. Wang and T. Gasser. Synchronizing sample curves nonparametrically. Ann. Statist., 27(2) :439-460, 1999.