

Foundations and Applications of Modern Nonparametric Statistics

Vladimir Spokoiny and Wolfgang Härdle

Weierstraß-Institute
for Applied Analysis and Stochastics
and Institute of Statistics and Economics,
Humboldt-Universität zu Berlin

October 4, 2006

Basics

Statistics is understanding data by modeling it.

Data $Y^{(n)} = (Y_1, \dots, Y_n)$ usually *random*.

$P = \mathcal{L}(Y^{(n)})$, the *unknown* joint distribution.

Statistical problem: to infer on P from the data $Y^{(n)}$.

Parametric modeling:

$$P = P_{\theta} \in (P_{\theta}, \theta \in \Theta \subset \mathbb{R}^p).$$

Nonparametric modeling: the parametric assumption is not fulfilled, or, equivalently, $p = \infty$.

Example

$$Y_1 = 1.2, Y_2 = 2.3, Y_3 = 0.6, \dots, Y_{20} = -1.5.$$

$$n = 20.$$

$$Y^{(n)} = (Y_1, \dots, Y_n).$$

Parametric modeling (for single observation): $P = P_\theta$, P_θ with pdf $\varphi(y - \theta)$.

Parametric modeling (for n observations): $P = P_\theta$, P_θ with pdf $\prod_{i=1}^n \varphi(y_i - \theta)$.

Example

$$\varphi(y) = (2\pi)^{-1/2} \exp(-y^2/2).$$

Shift model: P_{θ} is given by $\varphi(y - \theta)$, $p = 1$.

Shift scale model: P_{θ} is given by $\sigma^{-1}\varphi\{\sigma^{-1}(y - \mu)\}$,
 $\theta = (\mu, \sigma)$, $p = 2$.

Nonparametric modeling

$Y^{(n)} = (Y_1, \dots, Y_n)$ with pdf f from some function space, e.g.

$$f \in L_2 = \{f : \|f\|_2^2 < \infty, \int f = 1\}.$$

Fourier expansion:

$$f(y) = \sum_{j=1}^{\infty} \theta_j \psi_j(y),$$

$\{\psi_j\}$ ONB of the function space, $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots)$, $p = \infty$.

Parametric estimation. Maximum Likelihood approach

A parametric model: $Y^{(n)} \sim P_{\theta^*}$ for some $\theta^* \in \Theta$.

Let $P_{\theta} \ll P$ for some measure P for all $\theta \in \Theta$. Define the log-likelihood

$$L(\theta) = \log \frac{dP_{\theta}}{dP}(Y^{(n)}),$$

For some $\theta^{\circ} \in \Theta$, the (log)-likelihood ratio is

$$L(\theta, \theta^{\circ}) = L(\theta) - L(\theta^{\circ}) = \log \frac{dP_{\theta}}{dP_{\theta^{\circ}}}(Y^{(n)})$$

Maximum likelihood estimate (MLE) $\tilde{\theta}$ is the point of maximum of $L(\theta)$:

$$\tilde{\theta} = \operatorname{argmax}_{\theta \in \Theta} L(\theta) = \operatorname{argmax}_{\theta \in \Theta} L(\theta, \theta^{\circ}).$$

Minimum Contrast Estimate

Let $\mathcal{C}(\boldsymbol{\theta}) = \mathcal{C}(Y^{(n)}, \boldsymbol{\theta})$ be a *contrast function*.

$$\tilde{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \mathcal{C}(\boldsymbol{\theta}).$$

Example

The choice $\mathcal{C}(\boldsymbol{\theta}) = -L(\boldsymbol{\theta})$ leads back to the MLE.

Minimum Contrast Estimate. Example

P_θ = double exponential with pdf $\frac{1}{2} \exp(-|y|)$

Likelihood function $2^{-n} \prod_{i=1}^n \exp(-|Y_i - \theta|)$.

Log-likelihood $L(\theta) = -\sum_{i=1}^n |Y_i - \theta| - n \log 2$.

$$\begin{aligned} \mathcal{C}(\theta) &= \sum_{i=1}^n |Y_i - \theta| \quad \left[+n \log 2 \right] \\ L(\theta) &= -\mathcal{C}(\theta) \end{aligned}$$

Minimum Contrast Estimate. Example

Example

Least Squares Estimate (LSE) and Least Absolute Deviation (LAD) in the parametric regression model $E(Y_i|X_i) = f(X_i, \theta)$:

$$\tilde{\theta}_{LSE} = \operatorname{argmin}_{\theta \in \Theta} \sum_{i=1}^n \{Y_i - f(X_i, \theta)\}^2$$

$$\tilde{\theta}_{LAD} = \operatorname{argmin}_{\theta \in \Theta} \sum_{i=1}^n |Y_i - f(X_i, \theta)|$$

Some examples of parametric families. “Gaussian shift”

Let Y_1, \dots, Y_n be i.i.d. and follow

$$Y_i = \theta^* + \varepsilon_i$$

where $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ with known variance σ^2 .

Then

$$\begin{aligned} L(\theta) &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \theta)^2. \\ \tilde{\theta} &= \operatorname{argmax}_{\theta} L(\theta) = n^{-1} \sum_{i=1}^n Y_i = \theta^* + \sigma n^{-1/2} \xi, \end{aligned} \quad (1)$$

where ξ is standard normal.

“Gaussian shift”. Risk and confidence interval

The MLE $\tilde{\theta}$ fulfills $\tilde{\theta} = \theta^* + \sigma n^{-1/2} \xi$.

For any $r > 0$

$$\mathbf{E}_{\theta^*} |\tilde{\theta} - \theta^*|^r = \sigma^r n^{-r/2} c_r$$

with $c_r = \mathbf{E} |\xi|^r$. E.g. $c_1 = \sqrt{2/\pi}$, $c_2 = 1$ and $c_3 = 2c_1 = \sqrt{8/\pi}$, $c_4 = 3$.

If z_α is such that $P(|\xi| \leq z_\alpha) = 1 - \alpha$ with some $\alpha \in (0, 1)$, then

$$\mathcal{E}_\alpha = [\tilde{\theta} - \sigma n^{-1/2} z_\alpha, \tilde{\theta} + \sigma n^{-1/2} z_\alpha], \quad (2)$$

is a α -confidence interval for the parameter θ^* .

Confidence set

A **confidence set** (CS) to level α is a random set \mathcal{E}_α such that

$$P_{\theta^*}(\mathcal{E}_\alpha \ni \theta^*) \geq 1 - \alpha.$$

The \mathcal{E}_α from (2) is a CS to level α :

$$\begin{aligned} P_{\theta^*}(\mathcal{E}_\alpha \ni \theta^*) &= P_{\theta^*}(\tilde{\theta} - \sigma n^{-1/2} z_\alpha \leq \theta^* \leq \tilde{\theta} + \sigma n^{-1/2} z_\alpha) \\ &= P_{\theta^*}(\sqrt{n}|\tilde{\theta} - \theta^*|/\sigma \leq z_\alpha) = 1 - \alpha \end{aligned}$$

since (1) yields $\xi = \sqrt{n}|\tilde{\theta} - \theta^*|/\sigma$.

“Gaussian shift”. Fitted likelihood

For any θ

$$L(\tilde{\theta}, \theta) = L(\tilde{\theta}) - L(\theta) = n\sigma^{-2}(\tilde{\theta} - \theta)^2/2. \quad (3)$$

Since $\tilde{\theta} = \theta^* + \sigma n^{-1/2}\xi$,

$$2L(\tilde{\theta}, \theta^*) = n\sigma^{-2}(\tilde{\theta} - \theta^*)^2 = \xi^2 \sim \chi_1^2$$

If \mathfrak{z}_α is the α -quantile of χ_1^2 with $P(\xi^2 > \mathfrak{z}_\alpha) = \alpha$, then

$$\mathcal{E}_\alpha = \{u : 2L(\tilde{\theta}, u) \leq \mathfrak{z}_\alpha\} \quad (4)$$

is again an α -CS, but this time likelihood based.

Fitted likelihood “Gaussian shift”. Details

The formula $L(\tilde{\theta}, \theta) = n\sigma^{-2}(\tilde{\theta} - \theta)^2/2$ can be seen as follows:

$$\begin{aligned}L(\tilde{\theta}) - L(\theta) &= -(2\sigma^2)^{-1} \sum_{i=1}^n (Y_i - \tilde{\theta})^2 + (2\sigma^2)^{-1} \sum_{i=1}^n (Y_i - \theta)^2 \\ &= n\sigma^{-2}(\tilde{\theta} - \theta)^2/2 - \sigma^{-2} \sum_{i=1}^n (Y_i - \tilde{\theta})(\tilde{\theta} - \theta).\end{aligned}$$

The last sum is equal to zero by definition of $\tilde{\theta}$.

Bernoulli model

Let Y_1, \dots, Y_n be i.i.d. Bernoulli r.v.'s satisfying $P(Y_i = 1) = \theta$ and $P(Y_i = 0) = 1 - \theta$. Then

$$\begin{aligned}L(\theta) &= \log \prod_{i=1}^n \theta^{Y_i} (1 - \theta)^{1 - Y_i} = \log \theta \sum_i Y_i + \log(1 - \theta) \sum_i (1 - Y_i) \\ &= S \log \frac{\theta}{1 - \theta} + n \log(1 - \theta),\end{aligned}$$

where $S = Y_1 + \dots + Y_n$. Hence,

$$\tilde{\theta} = S/n, \quad \text{and} \quad L(\tilde{\theta}, \theta) = n\tilde{\theta} \log \frac{\tilde{\theta}}{\theta} + n(1 - \tilde{\theta}) \log \frac{1 - \tilde{\theta}}{1 - \theta} = n\mathcal{K}(\tilde{\theta}, \theta)$$

where $\mathcal{K}(\theta, \theta') = \theta \log(\theta/\theta') + (1 - \theta) \log\{(1 - \theta)/(1 - \theta')\}$ means the Kullback-Leibler divergence for the Bernoulli law.

Kullback-Leibler divergence

Kullback-Leibler (KL) divergence measures a “distance” between two distributions:

$$\mathcal{K}(P, Q) = \mathbf{E}_P \left\{ \log \left(\frac{dP}{dQ} \right) \right\}.$$

In terms of parametric model \mathbf{P}_θ :

$$\mathcal{K}(\theta, \theta') = \mathbf{E}_\theta \left\{ \log \left(\frac{d\mathbf{P}_\theta}{d\mathbf{P}_{\theta'}} \right) \right\}.$$

With pdf $p(y, \theta)$:

$$\mathcal{K}(\theta, \theta') = \mathbf{E}_\theta \left\{ \log \frac{p(y, \theta)}{p(y, \theta')} \right\} = \mathbf{E}_\theta \ell(\theta, \theta'), \quad \ell(\theta, \theta') = \log \frac{p(y, \theta)}{p(y, \theta')}.$$

Note that

$$\mathbf{E}_\theta L(\theta, \theta') = n \mathbf{E}_\theta \ell(\theta, \theta')$$

Poisson model

Let Y_1, \dots, Y_n be i.i.d. Poisson r.v.'s satisfying $P(Y_i = m) = \theta^m e^{-\theta} / m!$ for $m = 0, 1, 2, \dots$. Then

$$\begin{aligned} L(\theta) &= \log \prod_{i=1}^n \theta^{Y_i} e^{-\theta} / Y_i! = \log \theta \sum_{i=1}^n Y_i - n\theta - \log(Y_i!) \\ &= S \log \theta - n\theta + R, \end{aligned}$$

where $S = Y_1 + \dots + Y_n$ and $R = \sum_{i=1}^n \log(Y_i!)$. Therefore,

$$\tilde{\theta} = S/n, \quad \text{and} \quad L(\tilde{\theta}, \theta) = n\tilde{\theta} \log(\tilde{\theta}/\theta) - n(\tilde{\theta} - \theta) = n\mathcal{K}(\tilde{\theta}, \theta)$$

where $\mathcal{K}(\theta, \theta') = \theta \log(\theta/\theta') - (\theta - \theta')$ means the Kullback-Leibler divergence for the Poisson law.

Poisson model. Details

$$L(\tilde{\theta}, \theta) = n\tilde{\theta} \log(\tilde{\theta}/\theta) - n(\tilde{\theta} - \theta)$$

$$\begin{aligned}n\mathcal{K}(\theta, \theta') &= \mathbf{E}_{\theta} L(\theta, \theta') \\&= \mathbf{E}_{\theta} L(\theta) - \mathbf{E}_{\theta} L(\theta') \\&= \mathbf{E}_{\theta} [S \log \theta - n\theta] - \mathbf{E}_{\theta} [S \log \theta' - n\theta'] \\&= \mathbf{E}_{\theta} S \log(\theta/\theta') - n(\theta - \theta') \\&= n\theta \log(\theta/\theta') - n(\theta - \theta')\end{aligned}$$

Exponential model

Let Y_1, \dots, Y_n be i.i.d. exponential r.v.'s with parameter $\theta > 0$:
 $P(Y_i > t) = e^{-t/\theta}$. Then

$$L(\theta) = -n \log \theta - \sum_{i=1}^n Y_i / \theta = -S / \theta - n \log \theta,$$

where $S = Y_1 + \dots + Y_n$. Therefore

$$\tilde{\theta} = S/n, \quad \text{and} \quad L(\tilde{\theta}, \theta) = -n(1 - \tilde{\theta}/\theta) - n \log(\tilde{\theta}/\theta) = n\mathcal{K}(\tilde{\theta}, \theta)$$

where $\mathcal{K}(\theta, \theta') = \theta/\theta' - 1 - \log(\theta/\theta')$ means the Kullback-Leibler divergence for the exponential law.

Volatility model

Let ξ_1, \dots, ξ_n be i.i.d. $\mathcal{N}(0, \theta)$ r.v.'s and we observe $Y_i = \xi_i^2$.
Then

$$L(\theta) = -\frac{n}{2} \log(2\pi\theta) - \sum_{i=1}^n Y_i/(2\theta) = -\frac{n}{2} \log(2\pi\theta) - S/(2\theta),$$

where $S = Y_1 + \dots + Y_n$. Therefore,

$$\tilde{\theta} = S/n \quad \text{and} \quad L(\tilde{\theta}, \theta) = -\frac{n}{2} \log(\tilde{\theta}/\theta) - \frac{n}{2}(1 - \tilde{\theta}/\theta) = n\mathcal{K}(\tilde{\theta}, \theta)$$

where $\mathcal{K}(\theta, \theta') = 0.5(\theta/\theta' - 1) - 0.5 \log(\theta/\theta')$ means the Kullback-Leibler divergence for the two zero mean normal laws with the variance θ' and θ .

Exponential family

Means that all measures P_θ have density functions $p(y, \theta) = dP_\theta/dP(y)$ of the form

$$p(y, \theta) = p(y)e^{yC(\theta) - B(\theta)}.$$

Here $C(\theta)$ and $B(\theta)$ are some given nondecreasing functions on Θ and $p(y)$ is some nonnegative function on \mathcal{Y} .

The *natural* parametrization means the relation $E_\theta Y = \theta$. This and the identity $E p(y, \theta) \equiv 1$ yield

$$B'(\theta) = \theta C'(\theta).$$

Moreover, the Kullback-Leibler divergence

$\mathcal{K}(\theta, \theta') = E_\theta \log\{p(Y, \theta)/p(Y, \theta')\}$ for $\theta, \theta' \in \Theta$ satisfies

$$\mathcal{K}(\theta, \theta') = \theta\{C(\theta) - C(\theta')\} - \{B(\theta) - B(\theta')\}. \quad (5)$$

MLE for exponential family model

The density $p(y, \theta) = p(y)e^{yC(\theta) - B(\theta)}$ leads to the log-likelihood

$$L(\theta) = \sum_{i=1}^n \log p(Y_i, \theta) = SC(\theta) - nB(\theta) + R \quad (6)$$

where

$$S = \sum_{i=1}^n Y_i, \quad R = \sum_{i=1}^n \log p(Y_i).$$

The estimating equation $L'(\theta) = 0$ for $\tilde{\theta}$ and the identity $B'(\theta) \equiv \theta C'(\theta)$ result in

$$\tilde{\theta} = S/n = n^{-1} \sum_{i=1}^n Y_i \quad \text{and} \quad L(\tilde{\theta}, \theta) = L(\tilde{\theta}) - L(\theta) = n\mathcal{K}(\tilde{\theta}, \theta).$$

Exponential bound for the fitted likelihood. Exponential family

Theorem (Polzehl and Spokoiny (2005))

Let (P_θ) be an exponential family. Then for any $\mathfrak{z} > 0$

$$P_{\theta^*}(L(\tilde{\theta}, \theta^*) > \mathfrak{z}) \leq 2e^{-\mathfrak{z}}.$$

Moreover, for any $r > 0$

$$E_{\theta^*} L^r(\tilde{\theta}, \theta^*) = n^r E_{\theta^*} \mathcal{K}^r(\tilde{\theta}, \theta^*) \leq \mathfrak{r}_r,$$

where $\mathfrak{r}_r = 2r \int_{\mathfrak{z} \geq 0} \mathfrak{z}^{r-1} e^{-\mathfrak{z}} d\mathfrak{z} = 2r\Gamma(r)$.

θ^* maximizes $E_{\theta^*} L(\theta, \theta^*)$ over θ and $E_{\theta^*} L(\theta^*, \theta^*) = 0$.

$\tilde{\theta}$ maximizes $L(\theta, \theta^*)$ over θ , thus $L(\tilde{\theta}, \theta^*) \geq 0$. By Theorem 3 $L(\tilde{\theta}, \theta^*)$ is stochastically bounded.

A general exponential bound for the fitted likelihood

Theorem (Golubev and Spokoiny (2006))

Let $Y^{(n)}$ be from a parametric family $(P_{\theta}, \theta \in \Theta)$. Then under some regularity conditions for any $r > 0$

$$E_{\theta^*} L^r(\tilde{\theta}, \theta^*) \leq \mathfrak{R}_r,$$

where \mathfrak{R}_r depends on the parametric family only.

Moreover, if $\mathcal{V}^2(\theta, \theta') = \text{Var} L(\theta, \theta')$, then

$$E_{\theta^*} \mathcal{V}^{2r}(\tilde{\theta}, \theta^*) \leq \mathfrak{R}_r^*.$$

Some corollaries

In the regular case,

$$L(\tilde{\theta}, \theta^*) \approx \frac{n}{2} (\tilde{\theta} - \theta^*)^\top I(\theta^*) (\tilde{\theta} - \theta^*)$$

where $I(\theta^*)$, the Fisher information matrix, and the theorem yields **root-n consistency**:

$$E_{\theta^*}^{1/r} |\sqrt{I(\theta^*)}(\tilde{\theta} - \theta^*)|^r \leq C/\sqrt{n}.$$

Confidence sets: define

$$\mathcal{E}(\beta) = \{\theta : L(\tilde{\theta}, \theta^*) \leq \beta\}.$$

If β is sufficiently large then

$$P_{\theta^*}(\mathcal{E}(\beta) \ni \theta^*) \approx 0.$$

Advantages and problems of the parametric modeling

1. Well developed algorithms
2. Nice asymptotic theory. Root-n consistence and asymptotic normality of the estimator $\hat{\theta}$ under mild regularity assumptions.
3. Good in-sample properties.

Problem:

1. The parametric structure is crucial. If the parametric assumption is violated, the MLE estimator $\tilde{\theta}$ is often completely misspecified.

Aim: To extend the parametric approach and methods to the situation when the parametric assumption is not precisely fulfilled.