

# GAUSSIAN APPROXIMATION FOR NONPARAMETRIC MODELS

ION G. GRAMA

ABSTRACT. We show that the nonparametric model generated by independent observations  $X_i$ ,  $i = 1, \dots, n$ , with densities  $p(x, f(i/n))$ ,  $i = 1, \dots, n$ , where  $f$  is an "unknown" function  $f : [0, 1] \rightarrow \Theta$  in a smoothness class, can be approximated, in the sense of the Le Cam deficiency pseudodistance, by a nonparametric Gaussian shift model.

## 1. INTRODUCTION

Following Brown and Low [1], we say that two sequences of statistical experiments  $\mathcal{E}^n$ ,  $n = 1, 2, \dots$  and  $\mathcal{G}^n$ ,  $n = 1, 2, \dots$  are *asymptotically equivalent* if

$$\Delta(\mathcal{E}^n, \mathcal{G}^n) \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

where  $\Delta(\mathcal{E}^n, \mathcal{G}^n)$  is the Le Cam deficiency pseudo-distance between statistical experiments  $\mathcal{E}^n$  and  $\mathcal{G}^n$ . From the definition of the Le Cam distance (see Le Cam [5] or Le Cam and Yang [6]) it follows that asymptotically minimax risks over bounded loss functions in one sequence of models can be transferred to the asymptotically equivalent sequence of models, which means that one can compute the asymptotically minimax risk in non-Gaussian models by computing it in the accompanying Gaussian models. Thus asymptotic equivalence can be seen as an analog of the Hajék-Le Cam asymptotic minimax bound in the case of nonparametric models.

Such an asymptotic equivalence between the white noise model and its discrete time analog has been first established by Brown and Low [1]. For the density estimation, Nussbaum [7] found that the accompanying experiment is a white noise model, where the square root of the density is estimated. The more general case of generalized linear models was considered in Grama and Nussbaum [3]. However, in all these models it is assumed, essentially, that the noises are from an exponential family. The purpose of the present work is to extend the area of applicability of the asymptotic equivalence to the case of regression models with arbitrary noises.

## 2. THE RESULTS

**2.1. The model.** We start with a parametric model. Let  $\Theta$  be an interval (possibly infinite) in the real line and

$$\mathcal{E} = (X, \mathcal{X}, \{P_\theta : \theta \in \Theta\})$$

be a statistical experiment on the measurable space  $(X, \mathcal{X})$  with the parameter set  $\Theta$  and a dominating  $\sigma$ -finite measure  $\mu$ . For any  $\theta \in \Theta$ , denote by

$$p(x, \theta) = P_\theta(dx) / \mu(dx), \quad x \in X,$$

the Radon-Nikodym derivative of  $P_\theta$  w.r.t.  $\mu$ . For the sake of brevity we set  $p(\theta) = p(\cdot, \theta)$ . We shall assume in the sequel that  $p(\theta) > 0$ ,  $\mu$ -a.s. on  $X$ , for any  $\theta \in \Theta$ , which implies that the measures  $P_\theta$ ,  $\theta \in \Theta$  are equivalent:  $P_\theta \sim P_u$ , for  $\theta, u \in \Theta$ .

Now we shall introduce our nonparametric model. For any  $\beta > 0$ , let  $\Sigma^\beta$  be the Hölder ball of functions defined on  $[0, 1]$  and with values in  $\Theta$ , i.e. the set of functions  $f : [0, 1] \rightarrow \Theta$ , which satisfy Hölder's condition with exponent  $\beta$  and constant  $L$ :

$$|f^{(\beta_0)}(t) - f^{(\beta_0)}(s)| \leq L |t - s|^{\beta_1}, \quad t, s \in [0, 1] \quad \text{and} \quad \sup_{t \in T} |f(t)| \leq L,$$

where the integer  $\beta_0$  and the real  $0 < \beta_1 \leq 1$  are so that  $\beta = \beta_0 + \beta_1$ . The non-parametrically driven model, which we are going to deal with, assumes that we observe a sequence of independent r.v.'s  $X_1, \dots, X_n$ , with values in the measurable space  $(X, \mathcal{X})$ , so that, for each  $i = 1, \dots, n$ , the observation  $X_i$  has the density  $p(x, f(i/n))$ , where the function  $f$  is unknown and satisfies the smoothness condition  $f \in \Sigma^\beta$ . We shall make use of the notation  $P_f^n = P_{f(1/n)} \times \dots \times P_{f(n/n)}$ , where  $P_\theta$  is the probability measure in the parametric experiment  $\mathcal{E}$  and  $f \in \Sigma$ .

**2.2. Regularity assumptions.** Assume that  $\beta > 1/2$ . In the sequel we shall impose, on the density  $p(x, \theta)$  in the parametric experiment  $\mathcal{E}$ , the regularity assumptions (R1-R3), which are assumed to hold true with the same  $\varepsilon > 0$ .

**R1:** The function  $s(\theta) = \sqrt{p(\theta)}$  is smooth in the space  $L^2(X, \mathcal{X}, \mu)$ : there is a real number  $\delta \in (\frac{1}{2\beta}, 1)$  and a map  $\dot{s}(\theta) : \Theta \rightarrow L^2(X, \mathcal{X}, \mu)$  such that

$$\sup_{(\theta, u)} \frac{1}{|u - \theta|^{1+\delta}} \left( \int_X \left( s(x, u) - s(x, \theta) - (u - \theta) \dot{s}(x, \theta) \right)^2 \mu(dx) \right)^{1/2} < \infty,$$

where sup is taken over all pairs  $(\theta, u)$  satisfying  $\theta, u \in \Theta$ ,  $|u - \theta| \leq \varepsilon$ .

It is well-known (see Strasser [8]) that the function  $\dot{s}(\theta)$  in condition (R1) can be written as  $\dot{s}(\theta) = \frac{1}{2} \dot{l}(\theta) \sqrt{p(\theta)}$ ,  $\mu$ -a.s. on  $X$ , where  $\dot{l}(\theta) \in L^2(X, \mathcal{X}, \mu)$ . Moreover,  $\dot{l}(\theta) \in L^2(X, \mathcal{X}, P_\theta)$  and  $E_\theta \dot{l}(\theta) = 0$ ,  $\theta \in \Theta$ , where  $E_\theta$  is the

expectation under  $P_\theta$ . The map  $\dot{l}(\theta)$  is called *tangent vector* at  $\theta$ . For any  $\theta \in \Theta$ , define an *extended tangent vector*  $\dot{l}_\theta(u)$ ,  $u \in \Theta$ , by setting

$$\dot{l}_\theta(x, u) = \begin{cases} \dot{l}(x, \theta), & \text{if } u = \theta, \\ \frac{2}{u-\theta} \left( \sqrt{\frac{p(x, u)}{p(x, \theta)}} - 1 \right), & \text{if } u \neq \theta. \end{cases}$$

**R2:** There is a real number  $\delta \in (\frac{2\beta+1}{2\beta-1}, \infty)$  such that

$$\sup_{(\theta, u)} \int_X \left| \dot{l}_\theta(x, u) \right|^{2\delta} p(x, \theta) \mu(dx) < \infty,$$

where sup is taken over all pairs  $(\theta, u)$  satisfying  $\theta, u \in \Theta$ ,  $|u - \theta| \leq \varepsilon$ .

The Fisher information in the local experiment  $\mathcal{E}$  is

$$I(\theta) = \int_X \left( \dot{l}(x, \theta) \right)^2 p(x, \theta) \mu(dx), \quad \theta \in \Theta.$$

**R3:** There are two real numbers  $I_{\min}$  and  $I_{\max}$  such that

$$0 < I_{\min} \leq I(\theta) \leq I_{\max} < \infty, \quad \theta \in \Theta.$$

**2.3. Local result.** First we state a local Gaussian approximation. For any  $f \in \Sigma^\beta$ , denote by  $\Sigma_f^\beta(r)$  the neighborhood of  $f$ , shifted to the origin:

$$\Sigma_f^\beta(r) = \{h : |h| \leq r, f + h \in \Sigma^\beta\}.$$

Set

$$\gamma_n = c(\beta) \left( \frac{\log n}{n} \right)^{\frac{\beta}{2\beta+1}},$$

where  $c(\beta)$  is a constant depending on  $\beta$ . By definition the *local experiment*

$$\mathcal{E}_f^n = (X^n, \mathcal{X}^n, \{P_{f+h}^n : h \in \Sigma_f^\beta(\gamma_n)\})$$

is generated by the sequence of independent r.v.'s  $X_1, \dots, X_n$ , with values in the measurable space  $(X, \mathcal{X})$ , so that, for each  $i = 1, \dots, n$ , the observation  $X_i$  has the density  $p(x, g(i/n))$ , where  $g = f + h$ ,  $h \in \Sigma_f^\beta(\gamma_n)$ .

**Theorem 2.1.** *Let  $\beta > 1/2$  and  $I(\theta)$  be the Fisher information in the parametric experiment  $\mathcal{E}$ . Assume that the density  $p(x, \theta)$  satisfies the regularity conditions (R1 – R3). For any  $f \in \Sigma^\beta$ , let  $\mathcal{G}_f^n$  be the local Gaussian experiment, generated by the observations*

$$Y_i^n = h(i/n) + I(f(i/n))^{-1/2} \varepsilon_i, \quad i = 1, \dots, n,$$

with  $h \in \Sigma_f^\beta(\gamma_n)$ , where  $\varepsilon_1, \dots, \varepsilon_n$  is a sequence of i.i.d. standard normal r.v.'s. Then, uniformly in  $f \in \Sigma$ , the sequence of local experiments  $\mathcal{E}_f^n$ ,  $n = 1, 2, \dots$

is asymptotically equivalent to the sequence of local Gaussian experiments  $\mathcal{G}_f^n$ ,  $n = 1, 2, \dots$  :

$$\sup_{f \in \Sigma^\beta} \Delta(\mathcal{E}_f^n, \mathcal{G}_f^n) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

**2.4. Global result.** By definition the *global experiment*

$$\mathcal{E}^n = (X^n, \mathcal{X}^n, \{P_f^n : f \in \Sigma^\beta\})$$

is generated by the sequence of independent r.v.'s  $X_1, \dots, X_n$ , with values in the measurable space  $(X, \mathcal{X})$ , so that, for each  $i = 1, \dots, n$ , the observation  $X_i$  has the density  $p(x, f(i/n))$ , where  $f \in \Sigma^\beta$ . We shall use the following assumptions:

**G1:** For any  $\beta > \frac{1}{2}$ , there is an estimator  $\hat{f}_n : X^n \rightarrow \Sigma^\beta$ , such that

$$\sup_{f \in \Sigma^\beta} P\left(\left\|\hat{f}_n - f\right\|_\infty \geq c\bar{\gamma}_n\right) \rightarrow 0, \quad n \rightarrow \infty,$$

for any sequence  $\bar{\gamma}_n \geq 0$  satisfying  $\bar{\gamma}_n^{-1} = o(\gamma_n^{-1})$ .

**G2:** The Fisher information  $I(\theta) : \Theta \rightarrow (0, \infty)$  satisfies Hölder's condition with exponent  $\alpha \in (\frac{1}{2\beta}, 1)$ .

The main result of the paper is the following theorem, which states a global Gaussian approximation for the sequence of experiments  $\mathcal{E}^n$ ,  $n = 1, 2, \dots$  in the sense of the Le Cam distance.

**Theorem 2.2.** *Let  $\beta > 1/2$  and  $I(\theta)$  be the Fisher information in the parametric experiment  $\mathcal{E}$ . Assume that the density  $p(x, \theta)$  satisfies the regularity conditions (R1-R3) and that conditions (G1-G2) are satisfied. Let  $\mathcal{G}^n$  be the Gaussian experiment generated by the observations*

$$Y_i^n = \Gamma(f(i/n)) + \varepsilon_i, \quad i = 1, \dots, n,$$

with  $f \in \Sigma^\beta$ , where  $\varepsilon_1, \dots, \varepsilon_n$  is a sequence of i.i.d. standard normal r.v.'s and  $\Gamma(\theta) : \Theta \rightarrow R$  is any function satisfying  $\Gamma'(\theta) = \sqrt{I(\theta)}$ . Then the sequence of experiments  $\mathcal{E}^n$ ,  $n = 1, 2, \dots$  is asymptotically equivalent to the sequence of Gaussian experiments  $\mathcal{G}^n$ ,  $n = 1, 2, \dots$  :

$$\Delta(\mathcal{E}^n, \mathcal{G}^n) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

An example in Brown and Low [1] shows that asymptotic equivalence, in general, fails to hold true when  $\beta < \frac{1}{2}$ .

Assumption (G1) is of technical nature. It can be easily checked in the most particular cases of interest. We refer the reader to Section 3 for some examples.

The function  $\Gamma(\theta)$  can be related to so called *variance-stabilizing* transformation, which we proceed to introduce. Let  $X_1, \dots, X_n$  be a sequence of real valued i.i.d. r.v.'s, which depends on the parameter  $\theta \in \Theta$ . Let  $\theta$  be the common mean and  $\sigma(\theta)$  be the common variance. By the central limit theorem,

$$\sqrt{n}\{S_n - \theta\} \xrightarrow{d} N(0, \sigma(\theta)),$$

where  $S_n = (X_1 + \dots + X_n)/n$ . The variance-stabilizing transformation is defined to be a function  $F$  on the real line, which stabilize the variance in the limiting normal law, i.e. such that

$$\sqrt{n} \{F(S_n) - F(\theta)\} \xrightarrow{d} N(0, 1).$$

In this case the function  $\Gamma$  just coincides with the function  $F : \Gamma(\theta) = F(\theta)$ ,  $\theta \in \Theta$ . If the parametrization in the model  $\mathcal{E}$  is so that  $\mu'(\theta) = \sigma(\theta) = I(\theta)$ , where  $\mu(\theta) = E_\theta X$  is the mean and  $\sigma(\theta) = \text{Var}_\theta X$  of an observation  $X$  in the experiment  $\mathcal{E}$ , then  $\Gamma(\theta) = F(\mu(\theta))$ ,  $\theta \in \Theta$ .

Theorem 2.2 is proved using the local result in Theorem 2.1 by means of a globalizing procedure. The main idea of the proof of Theorem 2.1 is to decompose the initial experiment  $\mathcal{E}^n$  into a product of independent experiments and to show that each component can be "well approximated" by means of its Gaussian counterpart. For the last we develop a general approach according to which any experiment  $\mathcal{E}^n$ , satisfying a nonparametric analog of the local asymptotic quadraticity (LAQ) condition, can be constructed on the same measurable space with a Gaussian experiment  $\mathcal{G}^n$ , such that the Hellinger distance between the corresponding probability measures converges to 0 at a certain rate, as  $n \rightarrow \infty$ . The main tool in establishing this result is the functional analog of the Hungarian construction developed in Grama and Nussbaum [2]. Then we are able to check the nonparametric LAQ conditions for the model under consideration. Our approach is similar to that in Grama and Nussbaum [3] or Nussbaum [7].

### 3. EXAMPLES

**3.1. Linear regression model.** Consider the nonparametric regression model with non-Gaussian noise

$$(3.1) \quad X_i = f(i/n) + \xi_i, \quad i = 1, \dots, n,$$

where  $\xi_1, \dots, \xi_n$  are i.i.d. r.v.'s of means 0 and finite variances, with density  $p(x)$  on the real line,  $f \in \Sigma^\beta$  and  $\Sigma^\beta$  is a Hölder ball on  $[0, 1]$  with exponent  $\beta > \frac{1}{2}$ . This model is a particular case of the nonparametrically driven model, introduced in Section 2, when  $p(x, \theta) = p(x - \theta)$  is the shifted density  $p(x)$  and  $\theta \in R$ . It is easy to see that conditions (R1-R3) hold true, if we assume that the density  $p(x)$  obeys the following items:

**L1:** The function  $s(x) = \sqrt{p(x)}$  satisfy Hölder's condition with an exponent  $1 + \alpha$ , where  $\alpha \in (\frac{1}{2\beta}, 1)$ .

**L2:** For some  $\delta > \frac{2\beta+1}{2\beta-1}$  and  $\varepsilon > 0$ , we have

$$\sup_{|u| \leq \varepsilon} \int_{-\infty}^{\infty} |s'(x+u)/s(x)|^{2\delta} p(x) dx < \infty.$$

**L3:** The Fisher informational number is positive:

$$I = \int_{-\infty}^{\infty} (p'(x))^2/p(x)dx > 0.$$

It is well-known that a preliminary estimator satisfying condition (G1) exists. Then, under conditions (L1-L3), the nonparametric linear regression model (3.1) is asymptotically equivalent to a linear regression with Gaussian noise, in which we observe

$$Y_i = f(i/n) + I^{-1/2}\varepsilon_i, \quad i = 1, \dots, n,$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. standard normal r.v.'s.

**3.2. Exponential family model.** Other particular cases arise when the parametric experiment  $\mathcal{E} = (X, \mathcal{X}, \{P_\theta : \theta \in \Theta\})$  is an one-dimensional linearly indexed exponential family, where  $\Theta$  is a possibly infinite interval on the real line. The parameter set  $\Theta$  is assumed to be so that the regularity conditions (R1)-(R3) are satisfied. The parametrization in the examples below is so that an observation  $X$  in the experiment  $\mathcal{E}$  has the mean  $\mu(\theta) = \theta$ , since this form of the parametrization makes it easier to compute the function  $\Gamma(\theta)$ . Note that a preliminary estimator satisfying condition (G1) in the exponential family model exists (see Grama and Nussbaum [3]).

*Spectral density model.* Assume that we are given a sequence of normal observations  $X_1, \dots, X_n$  with means 0 and standard deviations  $f(i/n)$ , where the function  $f(t)$ ,  $t \in [0, 1]$  satisfies Hölder's condition with exponent  $\beta > \frac{1}{2}$  and is so that  $c_1 \leq f(t) \leq c_2$ , for some positive absolute constants  $c_1$  and  $c_2$ . In this model the density of the observations has the form  $p(x, \theta) = \frac{1}{\sqrt{2\pi\theta}} \exp\left(-\frac{x^2}{2\theta^2}\right)$ ,  $x \in R$  and the Fisher information is  $I(\theta) = 2\theta^{-2}$ . This gives us  $\Gamma(\theta) = \sqrt{2} \log \theta$ . Then the model is asymptotically equivalent to the Gaussian model, in which we observe

$$Y_i = \sqrt{2} \log f(i/n) + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. standard normal r.v.'s.

*Density model.* Assume that we are given a sequence of Poisson observations  $X_1, \dots, X_n$  with "unknown" intensities  $f(i/n)$ , where the function  $f(t)$ ,  $t \in [0, 1]$  satisfies Hölder's condition with exponent  $\beta > \frac{1}{2}$  and is so that  $c_1 \leq f(t) \leq c_2$ , for some positive absolute constants  $c_1$  and  $c_2$ . In this model  $p(x, \theta) = \theta^x \exp(-\theta)$ ,  $x \in X = \{0, 1, \dots\}$  and  $I(\theta) = \theta^{-1}$ . This gives us  $\Gamma(\theta) = 2\sqrt{\theta}$ . Then the model is asymptotically equivalent to the Gaussian model, in which we observe

$$Y_i = 2\sqrt{f(i/n)} + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. standard normal r.v.'s. We called this model the density model, since the density model reduces to the last one by using Katz poissonization technique (see Nussbaum [7]).

*Binary response model.* Assume that we are given a sequence of Bernoulli observations  $X_1, \dots, X_n$  taking values 0 and 1 with "unknown" probabilities of success  $f(i/n)$ , where the function  $f(t)$ ,  $t \in [0, 1]$  satisfies Hölder's condition with exponent  $\beta > \frac{1}{2}$  and is so that  $c_1 \leq f(t) \leq c_2$ , for some absolute constants  $c_1 > 0$  and  $c_2 < 1$ . In this model  $p(x, \theta) = \theta^x(1-\theta)^{1-x}$ ,  $x \in X = \{0, 1\}$  and  $I(\theta) = \frac{1}{\theta(1-\theta)}$ . This gives us  $\Gamma(\theta) = 2 \arcsin \sqrt{\theta}$ . Then this model is asymptotically equivalent to the Gaussian model, in which we observe

$$Y_i = 2 \arcsin \sqrt{f(i/n)} + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. standard normal r.v.'s.

#### REFERENCES

- [1] Brown, L. D. and Low, M. (1992). *Asymptotic equivalence of nonparametric regression and white noise*. Ann. Statist. **24**, 2384-2398.
- [2] Grama, I. G. and Nussbaum, M. (1997). *A nonstandard Hungarian construction for partial sums*. Preprint No. 324. Weierstrass-Institute, Berlin.
- [3] Grama, I. G. and Nussbaum, M. (1998). *Asymptotic equivalence for nonparametric generalized linear models*. Probability Theory and Related Fields (to appear). Preprint No. 289. Weierstrass-Institute, Berlin, 1986.
- [4] Ibragimov, I. A. and Has'minskii, R., H., (1981) *Statistical estimation. Asymptotic theory*. Springer-Verlag, New York etc.
- [5] Le Cam, L. (1986) *Asymptotic Methods in Statistical Decision Theory*. Springer-Verlag, New-York etc.
- [6] Le Cam, L. and Yang, G. (1990). *Asymptotics in Statistics: Some Basic Concepts*. Springer-Verlag, New-York etc.
- [7] Nussbaum, M. (1996) *Asymptotic equivalence of density estimation and white noise*. Ann. Statist. **24** 2399-2430.
- [8] Strasser, H. (1985) *Mathematical Theory of Statistics*. Walter de Gruyter, Berlin-New York.

Institute of Mathematics  
 Academy of Sciences  
 Academiei str. 5  
 277028 Kishinev  
 Moldova  
 e-mail: 16grama@mathem.moldova.su