

ON THE ASYMPTOTIC PROPERTIES OF A SIMPLE ESTIMATE OF THE MODE

CHRISTOPHE ABRAHAM¹, GÉRARD BIAU² AND BENOÎT CADRE³

Abstract. We consider an estimate of the mode θ of a multivariate probability density f with support in \mathbb{R}^d using a kernel estimate f_n drawn from a sample X_1, \dots, X_n . The estimate θ_n is defined as any x in $\{X_1, \dots, X_n\}$ such that $f_n(x) = \max_{i=1, \dots, n} f_n(X_i)$. It is shown that θ_n behaves asymptotically as any maximizer $\hat{\theta}_n$ of f_n . More precisely, we prove that for any sequence $(r_n)_{n \geq 1}$ of positive real numbers such that $r_n \rightarrow \infty$ and $r_n^d \log n/n \rightarrow 0$, one has $r_n \|\theta_n - \hat{\theta}_n\| \rightarrow 0$ in probability. The asymptotic normality of θ_n follows without further work.

Mathematics Subject Classification. 62G05.

Received May 15, 2002. Revised October 18, 2002.

INTRODUCTION

The problem of estimating the mode of a probability density has received considerable attention in the literature. For a historical and mathematical survey, we refer the reader to Sager [12]. One of the most recent application of mode estimation is in unsupervised *cluster analysis*, where one tries to break a complex data set into a series of piecewise similar groups or structures. The nonparametric approach is based on the premise that groups correspond to modes of a density. The goal then is to estimate the modes and assign each observation to the “domain of attraction” of a mode. But there are many other fields where the knowledge of the mode is of great interest. For example, the estimation of contours, or isopleths, is a natural extension of the estimation of modal points.

In this paper, we consider the problem of estimating the mode θ of a multivariate unimodal probability density f with support in \mathbb{R}^d from independent random variables X_1, \dots, X_n with density f . This problem has been studied by many authors, see for example Parzen [8], Konakov [5], Samanta [13], Devroye [2], Romano [10], Vieu [15], Leclerc and Pierre-Loti-Viaud [6], Mokkadem and Pelletier [7] and the references therein. Mostly, the estimate $\hat{\theta}_n$ of θ is defined as any maximizer of f_n , *i.e.*,

$$\hat{\theta}_n \in \operatorname{argmax}_{\mathbb{R}^d} f_n, \tag{0.1}$$

Keywords and phrases. Multivariate probability density, mode, kernel estimate, central limit theorem.

¹ ENSAM-INRA, UMR Biométrie et Analyse des Systèmes, 2 place Pierre Viala, 34060 Montpellier Cedex 1, France; e-mail: abraham@helios.ensam.inra.fr

² Laboratoire de Statistique Théorique et Appliquée, Université Pierre et Marie Curie – Paris VI, Boîte 158, 175 rue du Chevaleret, 75013 Paris, France; e-mail: biau@ccr.jussieu.fr

³ Laboratoire de Probabilités et Statistique, Université Montpellier II, Cc. 051, place Eugène Bataillon, 34095 Montpellier Cedex 5, France; e-mail: cadre@stat.math.univ-montp2.fr

where f_n is a *kernel density estimate* (Rosenblatt [11], Parzen [8], Devroye [3]). Recall that f_n is defined for all $x \in \mathbb{R}^d$ by

$$f_n(x) = \frac{1}{nh_n^d} \sum_{i=1}^n K\left(\frac{x - X_i}{h_n}\right),$$

where $(h_n)_{n \geq 1}$ is a sequence of positive real numbers such that $h_n \rightarrow 0$ and $K : \mathbb{R}^d \rightarrow \mathbb{R}$ is an integrable function with $\int_{\mathbb{R}^d} K(x) dx = 1$.

The estimate (0.1) is widely used, though it is hard to compute. Indeed, in addition to the calculation of f_n , it involves a numerical step for the computation of the argmax. As noticed by Devroye [2], classical search methods of the argmax perform satisfactorily only when f_n is sufficiently regular (continuous, unimodal, etc.) Thus, in practice, the argmax is usually computed over a finite grid. This failing is seldom discussed by authors, although it may affect the asymptotic properties of the estimate. Moreover, when the dimension of the sample space is large, or when accurate estimation is needed, the grid size (which exponentially increases with the dimension) leads to time-consuming computations. Finally, the search grid should be located around high density areas. In high dimension, this is a difficult task and the search grid usually includes low density areas.

As an attempt to remedy these problems, we proposed in a previous paper (Abraham *et al.* [1]) a concurrent estimate. Denoting by S_n the set $\{X_1, \dots, X_n\}$, we let the estimate θ_n be defined as

$$\theta_n \in \operatorname{argmax}_{S_n} f_n,$$

i.e.,

$$\theta_n \in \left\{ x \in S_n : f_n(x) = \max_{i=1, \dots, n} f_n(X_i) \right\}.$$

We emphasize that the main advantage of using θ_n instead of the argmax estimate (0.1) is that the former is easily computed in a finite number of operations. Moreover, since the sample points are naturally concentrated in high density areas, the set S_n can be regarded as the most natural (random) grid for approximating the mode. As pointed out by the referees, θ_n may also be an appropriate choice for a start value of any optimization algorithm to approximate $\hat{\theta}_n$. In [1], we established, under the condition $nh_n^d / \log n \rightarrow \infty$, the strong consistency of θ_n towards θ and provided an almost sure rate of convergence without any differentiability condition on f around the mode. This rate relies on the sharpness of the density near θ , which is measured by a *peak index*. For discussion, examples and numerical illustration, we refer the reader to [1].

One question still unanswered is whether the maximization over a finite sample alters the rate of convergence of the estimate θ_n compared to that of $\hat{\theta}_n$. In the present paper, we prove that the estimates θ_n and $\hat{\theta}_n$ have the same asymptotic behavior. In Section 1, we set up notation and assumptions and provide the main results. Proofs are gathered in Section 2.

1. NOTATION, HYPOTHESES AND MAIN RESULTS

1.1. Asymptotic proximity of θ_n and $\hat{\theta}_n$

Throughout the paper, we will denote by $\|\cdot\|$ the usual Euclidean norm for matrices or vectors and by $\operatorname{Hg}(x)$ the Hessian matrix at the point x of any function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ twice continuously differentiable in a neighborhood of x . The notation $\operatorname{diam} S$ will stand for the diameter of any set $S \subset \mathbb{R}^d$, *i.e.*,

$$\operatorname{diam} S = \sup_{x, y \in S} \|x - y\|.$$

For all $\varepsilon > 0$, the *level set* $A(\varepsilon)$ defined by

$$A(\varepsilon) = \{x \in \mathbb{R}^d : f(x) \geq f(\theta) - \varepsilon\}$$

will play a crucial role. If \xrightarrow{P} stands for the convergence in probability, we finally introduce the following hypotheses:

- H1** the application f is twice continuously differentiable on a neighborhood \mathcal{V} of θ and the matrix $Hf(\theta)$ is negative definite;
- H2** the convergence $\text{diam } A(\varepsilon) \rightarrow 0$ as $\varepsilon \rightarrow 0$ holds;
- H3** the sequence $(v_n(\hat{\theta}_n - \theta))_{n \geq 1}$ is tight for some sequence $(v_n)_{n \geq 1}$ of positive real numbers with $v_n \rightarrow \infty$;
- H4** the kernel K is twice continuously differentiable on \mathbb{R}^d and moreover $\sup_{x \in \mathcal{V}} \|Hf_n(x) - Hf(x)\| \xrightarrow{P} 0$;
- H5** one has $\sup_{x \in \mathbb{R}^d} |f_n(x) - \mathbf{E}f_n(x)| \xrightarrow{P} 0$.

Let us comment on these hypotheses. Assumption **H1** is a mild regularity assumption which is usually required to obtain rates of convergence in mode kernel estimation (see for example Parzen [8] and Romano [10]). Assumption **H2** has been introduced to avoid high density areas arbitrarily far from the mode. It can be shown that **H2** is equivalent to the classical condition

$$\sup_{x \notin \mathcal{U}} f(x) < f(\theta)$$

for any open vicinity \mathcal{U} of θ . For further discussion on this condition, we refer to Abraham *et al.* [1]. Assumption **H3** is a weak assumption which is in particular true when the sequence $(v_n(\hat{\theta}_n - \theta))_{n \geq 1}$ converges in distribution. In this respect, sufficient conditions are to be found in Romano [10] (for $d = 1$) and in Mokkadem and Pelletier [7] (for $d \geq 1$). Regarding **H4**, we refer to Silverman [14] for the univariate case and to Mokkadem and Pelletier [7] for the multivariate case. Assumption **H5** holds for example if K is of the form $K(x) = \Psi(\|x\|)$ where Ψ is a real valued function with bounded variation and $nh_n^d / \log n \rightarrow \infty$ (see Pollard [9], Th. 37, p. 34). More generally, it can be shown that **H5** holds whenever K satisfies a *covering number condition*, see Mokkadem and Pelletier [7] for a detailed discussion.

We are now ready to state the main result of the paper.

Theorem 1.1 (ASYMPTOTIC PROXIMITY OF θ_n AND $\hat{\theta}_n$). *Assume that **H1–H5** hold. For any sequence $(r_n)_{n \geq 1}$ of positive real numbers such that $r_n \rightarrow \infty$ and $r_n^d \log n / n \rightarrow 0$, we have*

$$r_n \|\theta_n - \hat{\theta}_n\| \xrightarrow{P} 0.$$

This theorem gains in interest if we realize that the very weak condition imposed on r_n allows to derive asymptotic properties of θ_n from analogous asymptotic properties of $\hat{\theta}_n$. Examples are presented in the next paragraph.

1.2. Application

The following corollary follows from Theorem 1.1 without further work. We let $\xrightarrow{\mathcal{D}}$ denote the convergence in distribution.

Corollary 1.1 (LIMIT LAW). *Assume that the assumptions **H1**, **H2**, **H4** and **H5** hold, $n^{d/2-1} h_n^{d(d+2)/2} \log n \rightarrow 0$ and*

$$\sqrt{nh_n^{d+2}} (\hat{\theta}_n - \theta) \xrightarrow{\mathcal{D}} Z,$$

where Z is some \mathbb{R}^d -valued random variable. Then

$$\sqrt{nh_n^{d+2}} (\theta_n - \theta) \xrightarrow{\mathcal{D}} Z.$$

The weak convergence of $\hat{\theta}_n$ to θ was first studied in the univariate framework by Parzen [8] who proved that if h_n is chosen such that $nh_n^6 \rightarrow \infty$ and $nh_n^7 \rightarrow 0$, then

$$\sqrt{nh_n^3}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{f(\theta)}{[f''(\theta)]^2} \int_{\mathbb{R}} K'^2(x) dx\right),$$

where \mathcal{N} is the Gaussian distribution. Eddy [4] and Romano [10] then proved that this central limit theorem still holds when the condition $nh_n^6 \rightarrow \infty$ is weakened to $nh_n^5/\log n \rightarrow \infty$. Recently, Mokkadem and Pelletier [7] extended these results to the multivariate framework. Precisely, these authors show, under suitable assumptions, that

$$\sqrt{nh_n^{d+2}}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, f(\theta)[\mathbf{H}f(\theta)]^{-1}\mathbf{G}[\mathbf{H}f(\theta)]^{-1}\right), \quad (1.1)$$

where \mathbf{G} is the $d \times d$ matrix defined by

$$G_{i,j} = \int_{\mathbb{R}^d} \frac{\partial K}{\partial x_i}(x) \frac{\partial K}{\partial x_j}(x) dx.$$

Therefore, under the assumptions of these authors, which imply **H1**, **H2**, **H4** and **H5**, the results above transfer to θ_n .

Following the remark of a referee, we would like to shed light on the fact that there are some problems associated with the use of results of this type. As an example, if one is interested in constructing confidence sets, it will be necessary to estimate the limiting variance matrix, which involves not only $f(\theta)$ but also the local sharpness around the peak, that is, the Hessian matrix $\mathbf{H}f(\theta)$. A possible answer is to use the weakly consistent estimates $f_n(\theta_n)$ and $\mathbf{H}f_n(\theta_n)$ of $f(\theta)$ and $\mathbf{H}f(\theta)$ as well as (1.1) in order to obtain, under suitable assumptions,

$$\sqrt{nh_n^{d+2}}\left[f_n(\theta_n)[\mathbf{H}f_n(\theta_n)]^{-1}\mathbf{G}[\mathbf{H}f_n(\theta_n)]^{-1}\right]^{-1/2}(\theta_n - \theta) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \mathbf{I}),$$

where \mathbf{I} denotes the $d \times d$ identity matrix.

2. PROOFS

2.1. Proof of Theorem 1.1

In the sequel, $B(a, \varepsilon)$ stands for the closed ball in $(\mathbb{R}^d, \|\cdot\|)$ with center at a and radius $\varepsilon > 0$. For all $n \geq 1$, $A_n(\varepsilon)$ will denote the random set

$$A_n(\varepsilon) = \left\{x \in \mathbb{R}^d : f_n(x) \geq f_n(\hat{\theta}_n) - \varepsilon\right\}.$$

First of all, we state two fundamental results. For the sake of clarity, their proofs are delayed to the end of the section.

Lemma 2.1. *Assume that **H1–H3** and **H5** hold. Then, for any sequence $(\alpha_n)_{n \geq 1}$ of positive random variables vanishing in probability, we have*

$$\text{diam } A_n(\alpha_n) \xrightarrow{P} 0.$$

Observe the correspondence between Lemma 2.1 and Assumption **H2**, where A and ε have been replaced by the random quantities A_n and α_n .

Proposition 2.1. *Assume that **H1** and **H3** hold, and let $(u_n)_{n \geq 1}$ be a sequence of positive real numbers such that $u_n \rightarrow 0$ and $nu_n^d/\log n \rightarrow \infty$. Then*

$$\mathbf{P}(\forall i \leq n : \|X_i - \hat{\theta}_n\| \geq u_n) \rightarrow 0.$$

The interest of Proposition 2.1 is in the assertion that there exists with high probability an observation within a distance u_n of $\hat{\theta}_n$. The main idea of the proof of Theorem 1.1 is to show that such an observation is also close to θ_n .

Proof of Theorem 1.1. Let $\eta > 0$ and, for all $n \geq 1$, $B_n = B(\hat{\theta}_n, c\eta/r_n)$, where $c \geq 1$ denotes a constant to be specified later. We have

$$\begin{aligned} \mathbf{P}(\theta_n \notin B_n) &\leq \mathbf{P}(\theta_n \notin B_n, \exists k \leq n : \|X_k - \hat{\theta}_n\| < \eta/r_n) \\ &\quad + \mathbf{P}(\forall k \leq n : \|X_k - \hat{\theta}_n\| \geq \eta/r_n). \end{aligned}$$

By Proposition 2.1, the second term of the right member vanishes as n grows to infinity. Thus, one only needs to prove that the first term tends to 0. Note first that, since $c \geq 1$, the event

$$[\theta_n \notin B_n, \exists k \leq n : \|X_k - \hat{\theta}_n\| < \eta/r_n]$$

is contained in the event

$$\left[\max_{X_i \notin B_n, i \leq n} f_n(X_i) \geq \max_{X_i \in B_n, i \leq n} f_n(X_i) \geq f_n(X_k) \right],$$

where the data X_k satisfies $\|X_k - \hat{\theta}_n\| < \eta/r_n$. Now, denoting by $\nabla f_n(x)$ the gradient of f_n at the point x , we have, according to the previous condition on X_k ,

$$\begin{aligned} |f_n(X_k) - f_n(\hat{\theta}_n)| &\leq \sup_{x \in B(\hat{\theta}_n, \eta/r_n)} \|\nabla f_n(x)\| \frac{\eta}{r_n} \\ &= \sup_{x \in B(\hat{\theta}_n, \eta/r_n)} \|\nabla f_n(x) - \nabla f_n(\hat{\theta}_n)\| \frac{\eta}{r_n} \\ &\quad (\text{since } \nabla f_n(\hat{\theta}_n) = 0) \\ &\leq a_n \left(\frac{\eta}{r_n} \right)^2, \end{aligned}$$

where $a_n = \sup_{x \in B(\hat{\theta}_n, \eta/r_n)} \|\mathbf{H}f_n(x)\|$. Consequently,

$$\begin{aligned} &\mathbf{P}(\theta_n \notin B_n, \exists k \leq n : \|X_k - \hat{\theta}_n\| < \eta/r_n) \\ &\leq \mathbf{P}\left(\max_{X_i \notin B_n, i \leq n} f_n(X_i) \geq f_n(\hat{\theta}_n) - a_n(\eta/r_n)^2 \right) \\ &\leq 1 - \mathbf{P}\left(\sup_{B_n^c} f_n < f_n(\hat{\theta}_n) - a_n(\eta/r_n)^2 \right). \end{aligned}$$

Hence, one only needs now to show that

$$\mathbf{P}\left(\sup_{B_n^c} f_n < f_n(\hat{\theta}_n) - a_n(\eta/r_n)^2 \right) \rightarrow 1 \quad \text{as } n \text{ tends to infinity.}$$

To this aim, observe that if $\alpha_n = 2a_n(\eta/r_n)^2$, the event

$$[\forall x \in B_n^c : f_n(x) < f_n(\hat{\theta}_n) - \alpha_n]$$

is contained in the event

$$\left[\sup_{B_n^c} f_n < f_n(\hat{\theta}_n) - a_n(\eta/r_n)^2 \right].$$

The first of the two events equals the event $[A_n(\alpha_n) \subset B_n]$. Consequently, the problem is reduced to showing that

$$\mathbf{P}\left(A_n(\alpha_n) \subset B_n\right) \rightarrow 1.$$

Using Taylor's formula and **H4** we have, for all $n \geq 1$ and $x \in A_n(\alpha_n)$,

$$f_n(x) - f_n(\hat{\theta}_n) = \frac{1}{2}(x - \hat{\theta}_n)^t \mathbf{R}_n(x)(x - \hat{\theta}_n), \quad (2.1)$$

where $\mathbf{R}_n(x) = \mathbf{H}f_n(x + \xi_x^n(x - \hat{\theta}_n))$ and $\xi_x^n \in (0, 1)$. Let us introduce the event \mathcal{E}_n defined by

$$\mathcal{E}_n = \left[\forall x \in A_n(\alpha_n) : (x - \hat{\theta}_n)^t \mathbf{R}_n(x)(x - \hat{\theta}_n) \leq \frac{1}{2}(x - \hat{\theta}_n)^t \mathbf{H}f(\theta)(x - \hat{\theta}_n) \right].$$

From (2.1), it is deduced that on the event \mathcal{E}_n ,

$$\forall x \in A_n(\alpha_n) : -(x - \hat{\theta}_n)^t \mathbf{H}f(\theta)(x - \hat{\theta}_n) \leq 4\alpha_n,$$

and consequently, that

$$A_n(\alpha_n) \subset B\left(\hat{\theta}_n, 2\sqrt{\alpha_n/\gamma}\right), \quad (2.2)$$

where $\gamma = \inf_{\|x\|=1} \|(-\mathbf{H}f(\theta))^{1/2}x\|^2$. Obviously $\gamma > 0$ according to **H1**. Therefore, according to (2.2), on the event \mathcal{E}_n , $A_n(\alpha_n) \subset B_n$ as soon as $c \geq 2\sqrt{2\alpha_n/\gamma}$. Thus, recalling that $c \geq 1$, we choose

$$c = 2 \max\left(1, 2\sqrt{\frac{2\|\mathbf{H}f(\theta)\|}{\gamma}}\right).$$

We obtain as sort

$$\begin{aligned} \mathbf{P}\left(A_n(\alpha_n) \subset B_n\right) &\geq \mathbf{P}\left(A_n(\alpha_n) \subset B_n, \mathcal{E}_n, c \geq 2\sqrt{2\alpha_n/\gamma}\right) \\ &= \mathbf{P}\left(\mathcal{E}_n, c \geq 2\sqrt{2\alpha_n/\gamma}\right). \end{aligned}$$

Since $\hat{\theta}_n \xrightarrow{P} \theta$ by **H3**, $\sup_{x \in \mathcal{V}} \|\mathbf{H}f_n(x) - \mathbf{H}f(x)\| \xrightarrow{P} 0$ by **H4** and $\mathbf{H}f$ is continuous on \mathcal{V} by **H1**, we get

$$\mathbf{P}\left(c \geq 2\sqrt{2\alpha_n/\gamma}\right) \rightarrow 1.$$

Thus, it remains to prove that $\mathbf{P}(\mathcal{E}_n) \rightarrow 1$. First note that

$$\begin{aligned} \sup_{x \in A_n(\alpha_n)} \|\mathbf{R}_n(x) - \mathbf{H}f(\theta)\| &\leq \sup_{x \in A_n(\alpha_n)} \left\| \mathbf{R}_n(x) - \mathbf{H}f(x + \xi_x^n(x - \hat{\theta}_n)) \right\| \\ &\quad + \sup_{x \in A_n(\alpha_n)} \left\| \mathbf{H}f(x + \xi_x^n(x - \hat{\theta}_n)) - \mathbf{H}f(\theta) \right\|. \end{aligned}$$

Observe now that

$$\sup_{x \in A_n(\alpha_n)} \left\| x + \xi_x^n(x - \hat{\theta}_n) - \hat{\theta}_n \right\| \leq 2 \text{diam } A_n(\alpha_n)$$

and that the bound vanishes in probability according to Lemma 2.1. Therefore, since $\hat{\theta}_n \xrightarrow{P} \theta$, we deduce from **H1** and **H4** that

$$\sup_{x \in A_n(\alpha_n)} \|\mathbf{R}_n(x) - \mathbf{H}f(\theta)\| \xrightarrow{P} 0. \quad (2.3)$$

Consequently,

$$\begin{aligned} \mathbf{P}(\mathcal{E}_n) &\geq \mathbf{P}(\forall x \in A_n(\alpha_n), \forall \|y\| = 1 : -y^t \mathbf{R}_n(x)y \geq -1/2 y^t \mathbf{H}f(\theta)y) \\ &\geq \mathbf{P}(\forall \|y\| = 1 : \sup_{x \in A_n(\alpha_n)} |y^t (\mathbf{R}_n(x) - \mathbf{H}f(\theta))y| \leq -1/2 y^t \mathbf{H}f(\theta)y) \\ &\quad (\text{using the triangle inequality}) \\ &\geq \mathbf{P}\left(\sup_{x \in A_n(\alpha_n)} \|\mathbf{R}_n(x) - \mathbf{H}f(\theta)\| \leq -1/2 \inf_{\|y\|=1} y^t \mathbf{H}f(\theta)y\right) \end{aligned} \quad (2.4)$$

where, for the last inequality, we used the fact that

$$\sup_{\|y\|=1} \sup_{x \in A_n(\alpha_n)} |y^t (\mathbf{R}_n(x) - \mathbf{H}f(\theta))y| \leq \sup_{x \in A_n(\alpha_n)} \|\mathbf{R}_n(x) - \mathbf{H}f(\theta)\|.$$

The probability (2.4) tends to 1 according to (2.3), since $\mathbf{H}f(\theta)$ is negative definite. Consequently, $\mathbf{P}(\mathcal{E}_n) \rightarrow 1$, hence the theorem is proved. \square

2.2. Proof of Lemma 2.1

For all $n \geq 1$ and $\mu > 0$, we set

$$\beta_n = \alpha_n + \sup_{x \in \mathbb{R}^d} |f_n(x) - \mathbf{E}f_n(x)| + |f(\theta) - f_n(\hat{\theta}_n)|$$

and

$$D_n(\mu) = \{x \in \mathbb{R}^d : \mathbf{E}f_n(x) \geq f(\theta) - \mu\}.$$

We note that the first two terms in β_n go to zero in virtue of **H5**. With respect to the third term, it is bounded by

$$|f(\theta) - f(\hat{\theta}_n)| + |f(\hat{\theta}_n) - f_n(\hat{\theta}_n)|.$$

The first of the two terms above tends to 0 in probability by the continuity of f around θ and the fact that $\hat{\theta}_n \xrightarrow{P} \theta$. Finally, the second term vanishes under **H5** and the uniform continuity of f around its mode.

Moreover, attention shows that

$$A_n(\alpha_n) \subset D_n(\beta_n).$$

Let $\varepsilon > 0$. Then, for all $\gamma > 0$, the following chain of inequalities is valid.

$$\begin{aligned} \mathbf{P}(\text{diam } A_n(\alpha_n) > 4\varepsilon) &\leq \mathbf{P}(\text{diam } D_n(\beta_n) > 4\varepsilon) \\ &\leq \mathbf{P}(\text{diam } D_n(\beta_n) > 4\varepsilon, \beta_n \leq \gamma) + \mathbf{P}(\beta_n \geq \gamma) \\ &\leq \mathbf{P}(\text{diam } D_n(\gamma) > 4\varepsilon) + \mathbf{P}(\beta_n \geq \gamma). \end{aligned}$$

As mentioned above, we have $\beta_n \xrightarrow{P} 0$. Thus, one only needs to prove the existence of $\gamma > 0$ such that, for all n large enough,

$$D_n(\gamma) \subset B(\theta, 2\varepsilon).$$

It is easy to deduce from **H2** (see also Abraham *et al.* [1], Lem. 1) the existence of $\gamma > 0$ such that

$$\sup_{B(\theta, \varepsilon)^c} f < f(\theta) - 2\gamma.$$

Now, since K is integrable, there exists a compact set $T \subset \mathbb{R}^d$ with

$$f(\theta) \int_{T^c} |K(y)| \, dy < \gamma.$$

Consequently, if n is large enough, we have, for all $x \in B(\theta, 2\varepsilon)^c$,

$$\begin{aligned} \mathbf{E}f_n(x) &= \int_T K(y)f(x - h_n y) \, dy + \int_{T^c} K(y)f(x - h_n y) \, dy \\ &\leq \sup_{B(\theta, \varepsilon)^c} f + f(\theta) \int_{T^c} |K(y)| \, dy \\ &< f(\theta) - \gamma. \end{aligned}$$

Hence, for all n large enough, $D_n(\gamma) \subset B(\theta, 2\varepsilon)$, and the proof of the lemma is complete. \square

2.3. Proof of Proposition 2.1

We have divided the proof of Proposition 2.1 into a sequence of two lemmas. Before stating these two lemmas, we need to introduce some additional notations. From now on, ρ denotes a fixed positive real number such that $\inf_{B(\theta, \rho)} f > 0$. Note that such a ρ does exist under Assumption **H1**. For any \mathbb{R}^d -valued random variable Z , we let Z^* be the random variable defined as follows:

$$Z^* = \begin{cases} \theta + \rho(Z - \theta)/\|Z - \theta\| & \text{if } Z \notin B(\theta, \rho) \\ Z & \text{if } Z \in B(\theta, \rho). \end{cases}$$

It is worth pointing out that Z^* is a *truncated version* of Z . Indeed, for $Z \notin B(\theta, \rho)$, Z^* is defined as the intersection of the line θZ with the sphere $\{x : \|x - \theta\| = \rho\}$.

We also introduce a *ghost sample* Y_1, \dots, Y_n of independent and identically distributed random variables with density f . This sample is assumed to be *independent* of the sample X_1, \dots, X_n . Finally, for any subsets $E, F \subset \mathbb{R}^d$, we define $\delta(E, F)$ as

$$\delta(E, F) = \sup_{x \in E} \inf_{y \in F} \|x - y\|.$$

The quantity $\delta(E, F)$ may be regarded as a distance between the sets E and F reminiscent of the usual Hausdorff metric.

Lemma 2.2. *Assume that **H1** holds and let $(u_n)_{n \geq 1}$ be a sequence of positive real numbers such that $u_n \rightarrow 0$ and $nu_n^d / \log n \rightarrow \infty$. Then*

$$\mathbf{P}\left(\delta(\{Y_1^*, \dots, Y_n^*\}, \{X_1^*, \dots, X_n^*\}) \geq u_n\right) \rightarrow 0.$$

Proof of Lemma 2.2. Set $Y = Y_1$ and $\delta_n^* = \delta(\{Y_1^*, \dots, Y_n^*\}, \{X_1^*, \dots, X_n^*\})$. We have, for all $n \geq 1$,

$$\begin{aligned}
\mathbf{P}(\delta_n^* \geq u_n) &= \mathbf{P}(\exists i \leq n, \forall j \leq n : \|Y_i^* - X_j^*\| \geq u_n) \\
&= 1 - \mathbf{P}(\forall i \leq n, \exists j \leq n : \|Y_i^* - X_j^*\| < u_n) \\
&= 1 - \mathbf{E} \left[\mathbf{P}(\exists j \leq n : \|Y^* - X_j^*\| < u_n | X_1^*, \dots, X_n^*)^n \right] \\
&\leq 1 - \mathbf{P}(\exists j \leq n : \|Y^* - X_j^*\| < u_n)^n \\
&\quad \text{(using Jensen's inequality)} \\
&= 1 - \left(1 - \mathbf{P}(\forall j \leq n : \|Y^* - X_j^*\| \geq u_n) \right)^n \\
&= 1 - \left(1 - \mathbf{E} \left[\mathbf{P}(\|Y^* - X^*\| \geq u_n | Y^*)^n \right] \right)^n,
\end{aligned}$$

where $X = X_1$. Hence, one only needs to prove that

$$n \mathbf{E} \left[\mathbf{P}(\|Y^* - X^*\| \geq u_n | Y^*)^n \right] \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

By the very definition of X^* , we have

$$\begin{aligned}
\mathbf{P}(\|Y^* - X^*\| \geq u_n | Y^*) &\leq 1 - \mathbf{P}(X^* \in B(Y^*, u_n/2) | Y^*) \\
&\leq 1 - \mathbf{P}(X^* \in B(Y^*, u_n/2), X \in B(\theta, \rho) | Y^*) \\
&= 1 - \mathbf{P}(X \in B(Y^*, u_n/2) \cap B(\theta, \rho) | Y^*) \\
&= 1 - \int_{B(Y^*, u_n/2) \cap B(\theta, \rho)} f(x) \, dx \\
&\leq 1 - \lambda(B(Y^*, u_n/2) \cap B(\theta, \rho)) \inf_{B(\theta, \rho)} f,
\end{aligned}$$

where λ denotes the Lebesgue measure on \mathbb{R}^d . Since $Y^* \in B(\theta, \rho)$, we deduce that for some constant $c > 0$ depending only on ρ, f and d , we have, with probability one,

$$\mathbf{P}(\|Y^* - X^*\| \geq u_n | Y^*) \leq 1 - c u_n^d.$$

Therefore

$$n \mathbf{E} \left[\mathbf{P}(\|Y^* - X^*\| \geq u_n | Y^*)^n \right] \leq n(1 - c u_n^d)^n,$$

and the bound vanishes under the conditions $u_n \rightarrow 0$ and $n u_n^d / \log n \rightarrow \infty$. \square

Lemma 2.3. *Assume that **H1** and **H3** hold, and let $(u_n)_{n \geq 1}$ be a sequence of positive real numbers such that $u_n \rightarrow 0$ and $n u_n^d \rightarrow \infty$. Then*

$$\mathbf{P}(\forall i \leq n : \|Y_i - \hat{\theta}_n\| \geq u_n) \rightarrow 0.$$

Proof of Lemma 2.3. Let $\varepsilon > 0$. By **H3**, there exists $a > 0$ such that $\sup_{n \geq 1} \mathbf{P}(v_n \|\hat{\theta}_n - \theta\| > a) \leq \varepsilon$. If, for all $n \geq 1$, T_n denotes the ball $B(\theta, a/v_n)$, ν_n denotes the law of $\hat{\theta}_n$ and $Y = Y_1$, we deduce from the independence of $\hat{\theta}_n$ and (Y_1, \dots, Y_n) that

$$\begin{aligned} \mathbf{P}(\forall i \leq n : \|Y_i - \hat{\theta}_n\| \geq u_n) &= \int \mathbf{P}(\forall i \leq n : \|Y_i - t\| \geq u_n) \nu_n(dt) \\ &\leq \varepsilon + \int_{T_n} \mathbf{P}(\|Y - t\| \geq u_n)^n \nu_n(dt) \\ &\leq \varepsilon + \sup_{t \in T_n} \mathbf{P}(\|Y - t\| \geq u_n)^n. \end{aligned}$$

Now, for all n large enough,

$$\begin{aligned} \sup_{t \in T_n} \mathbf{P}(\|Y - t\| \geq u_n)^n &= \sup_{t \in T_n} \left(1 - \int_{B(t, u_n)} f(x) dx\right)^n \\ &\quad \text{(using the compactness of } T_n) \\ &\leq \left(1 - c u_n^d \inf_{B(\theta, \rho)} f\right)^n, \end{aligned}$$

where $c > 0$ is a constant which only depends on the dimension d . Since $u_n \rightarrow 0$, $n u_n^d \rightarrow \infty$ and $\inf_{B(\theta, \rho)} f > 0$, the lemma is proved. \square

We are now in a position to prove Proposition 2.1.

Proof of Proposition 2.1. As in the proof of Lemma 2.2, we use the notation $\delta_n^* = \delta(\{Y_1^*, \dots, Y_n^*\}, \{X_1^*, \dots, X_n^*\})$. For all $n \geq 1$, we have

$$\begin{aligned} \mathbf{P}(\forall i \leq n : \|X_i - \hat{\theta}_n\| \geq u_n) &\leq \mathbf{P}(\forall i \leq n : \|X_i - \hat{\theta}_n\| \geq u_n, \hat{\theta}_n \in B(\theta, \rho/2)) \\ &\quad + \mathbf{P}(\hat{\theta}_n \notin B(\theta, \rho/2)). \end{aligned}$$

By **H3**, the last term vanishes. Moreover, assuming that n is large enough to ensure that $u_n \leq \rho/2$, we can write

$$\begin{aligned} &\mathbf{P}(\forall i \leq n : \|X_i - \hat{\theta}_n\| \geq u_n, \hat{\theta}_n \in B(\theta, \rho/2)) \\ &= \mathbf{P}(\forall i \leq n : \|X_i^* - \hat{\theta}_n\| \geq u_n, \hat{\theta}_n \in B(\theta, \rho/2)) \\ &\leq \mathbf{P}(\forall i \leq n : \|X_i^* - \hat{\theta}_n\| \geq u_n, \delta_n^* \leq u_n/2, \hat{\theta}_n \in B(\theta, \rho/2)) \\ &\quad + \mathbf{P}(\delta_n^* > u_n/2). \end{aligned}$$

By Lemma 2.2 the latter term tends to 0. Finally, by the very definition of the Y_j^* 's, we obtain

$$\begin{aligned}
& \mathbf{P}\left(\forall i \leq n : \|X_i^* - \hat{\theta}_n\| \geq u_n, \delta_n^* \leq u_n/2, \hat{\theta}_n \in B(\theta, \rho/2)\right) \\
&= \mathbf{P}\left(\forall i \leq n : \|X_i^* - \hat{\theta}_n\| \geq u_n, \forall j \leq n, \exists k \leq n : \|Y_j^* - X_k^*\| \leq u_n/2, \right. \\
&\quad \left. \hat{\theta}_n \in B(\theta, \rho/2)\right) \\
&\leq \mathbf{P}\left(\forall j \leq n, \exists k \leq n : \|X_k^* - \hat{\theta}_n\| \geq u_n, \|Y_j^* - X_k^*\| \leq u_n/2, \hat{\theta}_n \in B(\theta, \rho/2)\right) \\
&\leq \mathbf{P}\left(\forall j \leq n : \|Y_j^* - \hat{\theta}_n\| \geq u_n/2, \hat{\theta}_n \in B(\theta, \rho/2)\right) \\
&= \mathbf{P}\left(\forall j \leq n : \|Y_j - \hat{\theta}_n\| \geq u_n/2, \hat{\theta}_n \in B(\theta, \rho/2)\right),
\end{aligned}$$

where the last equality holds if n is large enough, to ensure that $u_n \leq \rho$. The lemma is then a straightforward consequence of Lemma 2.3. \square

Acknowledgements. The authors greatly thank the anonymous referees and an Associate Editor for a careful reading of the paper and many helpful comments.

REFERENCES

- [1] C. Abraham, G. Biau and B. Cadre, Simple estimation of the mode of a multivariate density. *Canadian J. Statist.* **31** (2003) 23-34.
- [2] L. Devroye, Recursive estimation of the mode of a multivariate density. *Canadian J. Statist.* **7** (1979) 159-167.
- [3] L. Devroye, *A Course in Density Estimation*. Birkhäuser, Boston (1987).
- [4] W.F. Eddy, Optimum kernel estimates of the mode. *Ann. Statist.* **8** (1980) 870-882.
- [5] V.D. Konakov, On asymptotic normality of the sample mode of multivariate distributions. *Theory Probab. Appl.* **18** (1973) 836-842.
- [6] J. Leclerc and D. Pierre-Loti-Viaud, Vitesse de convergence presque sûre de l'estimateur à noyau du mode. *C. R. Acad. Sci. Paris* **331** (2000) 637-640.
- [7] A. Mokkadem and M. Pelletier, A law of the iterated logarithm for the kernel mode estimator, *ESAIM: Probab. Statist.* **7** (2003) 1-21.
- [8] E. Parzen, On estimation of a probability density function and mode. *Ann. Math. Statist.* **33** (1962) 1065-1076.
- [9] D. Pollard, *Convergence of Stochastic Processes*. Springer-Verlag, New York (1984).
- [10] J.P. Romano, On weak convergence and optimality of kernel density estimates of the mode. *Ann. Statist.* **16** (1988) 629-647.
- [11] M. Rosenblatt, Remarks on some nonparametric estimates of a density function. *Ann. Math. Statist.* **27** (1956) 832-837.
- [12] T.W. Sager, Estimating modes and isopleths. *Comm. Statist. - Theory Methods* **12** (1983) 529-557.
- [13] M. Samanta, Nonparametric estimation of the mode of a multivariate density. *South African Statist. J.* **7** (1973) 109-117.
- [14] B. Silverman, Weak and strong uniform consistency of the kernel estimate of a density and its derivatives. *Ann. Statist.* **6** (1978) 177-184.
- [15] P. Vieu, A note on density mode estimation. *Statist. Probab. Lett.* **26** (1996) 297-307.