## Point process models and local asymptotics in statistics

## II – Local asymptotic normality – comparison of estimators

Reinhard Höpfner, University of Mainz

NOMP II, March 22–24, 2021

We explain local asymptotic normality (LAN) in the sense of Le Cam, and give the main theorems on comparison of estimators under LAN. I will not give the proofs. All results below admit an extension to the –important and more general– setting of local asymptotic mixed normality (LAMN), but I will give only some very few hints in this direction. Key references are

- Le Cam, L.: Théorie asymptotique de la décision statistique. Montreal 1969.
- Hájek, J.: A characterization of limiting distributions for regular estimators.
   Zeitschrift f. Wahrscheinlichkeitstheorie u. verw. Geb. 14, 232–330 (1970)
- Le Cam, L.: Limits of experiments. Proc. 6th Berkeley symposium on mathematical statistics and probability. Vol. I, 245–261. Univ. Calif. Press 1972.
- Davies, R.: Asymptotic inference when the amount of information is random.
  In: LeCam, L., Olshen, R., Eds: Proc. of the Berkeley Symposium in in honor of J. Neyman and J. Kiefer, Vol. II: Wadsworth, 1985.
- Le Cam, L., Yang, G.: Asymptotics in statistics: some basic concepts. Springer 1990.

See also the books

- Hájek, J., Sidak: Theory of rank tests. Academic Press 1967.
- Strasser, H.: Mathematical theory of statistics. deGruyter 1985.
- Pfanzagl, J.: Parametric statistical theory. deGruyter 1994
- Liese, F., Miescke, K.: Statistical decision theory. Springer 2008.
- Höpfner, R.: Asymptotic statistics with a view to stochastic processes. deGruyter 2014.

When I give precise links (to proofs, background results, or LAMN) I will refer to the last title.

Quoting the work of Le Cam and of Hájek up to  $\approx$  1975, the russian school took a different road to a local asymptotic minimax theorem:

- Ibragimov, I., Khasminskii, R.: Statistical estimation. Springer 1981.
- Kutoyants, Y.: Statistical inference in ergodic diffusion processes. Springer 2004.

# 6 L<sup>2</sup>-differentiable statistical models and iid experiments

This section deals with iid models. Consider dominated experiments

$$\mathcal{E} := (\Omega, \mathcal{A}, \mathcal{P} := \{P_{\vartheta} : \vartheta \in \Theta\}) \quad , \quad \mathcal{P} \ll \nu \,, \ f_{\vartheta} := \frac{dP_{\vartheta}}{d\nu} \,, \ \Theta \subset \mathbb{R}^d \text{ open}$$

and n-fold product experiments

$$\mathcal{E}_n := \left( \Omega_n := \bigwedge_{j=1}^n \Omega, \ \mathcal{A}_n := \bigotimes_{j=1}^n \mathcal{A}, \ \mathcal{P}_n := \left\{ P_{\vartheta}^n := \bigotimes_{j=1}^n P_{\vartheta} : \vartheta \in \Theta \right\} \right)$$

and let  $X = (X_1, \ldots, X_n)$  denote the canonical statistics on  $(\Omega_n, \mathcal{A}_n)$ .

# **Definition :** $\mathcal{E}$ is called <u>L<sup>2</sup>-differentiable at $\vartheta$ with derivative $V_{\vartheta}$ if there is some</u>

 $\mathbb{R}^d$ -valued random variable  $V_{\vartheta}$  with components  $V_{\vartheta,1}, \ldots, V_{\vartheta,d}$  in  $L^2(\Omega, \mathcal{A}, P_{\vartheta})$ 

such that

$$\frac{1}{|\xi - \vartheta|^2} \int_{\Omega} \left| f_{\xi}^{1/2} - f_{\vartheta}^{1/2} - \frac{1}{2} f_{\vartheta}^{1/2} \left( \xi - \vartheta \right)^{\top} V_{\vartheta} \right|^2 d\nu \quad \longrightarrow \quad 0 \quad \text{as} \ \xi \to \vartheta$$

We call

$$J(\vartheta) := E_{\vartheta} \left( V_{\vartheta}^{\top} V_{\vartheta} \right)$$

<u>Fisher-Information in  $\vartheta$ </u> (when  $L^2$ -differentiability holds,  $V_\vartheta$  does not depend on the choice of the dominating measure, and is centred:  $E_\vartheta(V_\vartheta) = 0$  in all components).

**Example:** i) Consider a location model  $\mathcal{E}$  on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  with Lebesgue densities

$$f_{\vartheta}(x) = \frac{1}{2}e^{-|x-\vartheta|} , \quad \vartheta \in \Theta := \mathbb{R} , \quad x \in \mathbb{R}$$

 $\mathcal{E}$  is  $L^2$ -differentiable at every  $\vartheta$  with derivative  $V_{\vartheta}(x) := \operatorname{sign}(x - \vartheta)$ . The Fisher information is  $J(\vartheta) = E_{\vartheta}(V_{\vartheta}^2) = 1$  for all  $\vartheta \in \Theta$ .

ii) In smoothly parametrized models we can prove  $L^2$ -differentiability under mild assumptions, with derivative  $V_{\vartheta}$  of type  $\frac{\partial}{\partial \vartheta} \log f_{\vartheta}(\cdot)$ . '2nd Le Cam lemma': Consider a point  $\vartheta$  where  $\mathcal{E}$  is  $L^2$ -differentiable with derivative  $V_\vartheta$ . Then for bounded sequences  $(h_n)$  in  $\mathbb{R}^d$  (we always assume that  $\vartheta + n^{-1/2}h_n$  belongs to  $\Theta$ ), log-likelihood ratios in  $\mathcal{E}_n$  admit at  $\vartheta$  quadratic expansions

$$\log L_n^{(\vartheta+n^{-1/2}h_n)/\vartheta} = h_n^{\top} S_n(\vartheta) - \frac{1}{2} h_n^{\top} J(\vartheta) h_n + o_{P_{\vartheta}^n}(1) \quad , \quad n \to \infty$$

with

$$S_n(\vartheta) := n^{-1/2} \sum_{j=1}^n V_{\vartheta}(X_j)$$

(called score at  $\vartheta$  in  $\mathcal{E}_n$ ) converging in law:

$$\mathcal{L}(S_n(\vartheta) \mid P_\vartheta^n) \xrightarrow{w} \mathcal{N}(0, J(\vartheta)) \quad (\text{weak convergence in } \mathbb{R}^d \text{ as } n \to \infty)$$

**Remark:** Compare this to normal distributions  $\mathcal{P} := \{ P_h := \mathcal{N}(Jh, J) : h \in \mathbb{R}^d \}$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d)), J \in \mathbb{R}^{d \times d}$  symmetric and strictly positive definite, where we have

$$\frac{dP_h}{dP_0}(x) = \exp\left\{h^{\top}S(x) - \frac{1}{2}h^{\top}Jh\right\} , \quad S(x) := x$$

with  $\mathcal{L}(S \mid P_0) = \mathcal{N}(0, J).$ 

#### Proofs in many books, e.g. sections 4.1 and 4.2 in H. 2014.

Results of type '2nd Le Cam Lemma' can be proved in broad classes of stochastic process models under ergodicity assumptions: then martingales and their angle brackets appear in place of score and Fisher information, and one applies martingale convergence theorems (e.g. Jacod-Shiryaev 1987, VIII.3.22).

- Löcherbach, E.: LAN and LAMN for systems of interacting diffusions with branching and immigration. Ann. H. Poincaré Proba. Stat. 38, 59–90 (2002).
- Holbach, S.: Local asymptotic normality for shape and periodicity of a signal in the drift of a degenerate diffusion with internal variables. Electronic J. Probability **13**, 4884–4915 (2019).

In a variety of other stochastic process models, results of type '2nd Le Cam Lemma' exist, with information processes (to which limit theorems apply) in place of deterministic information.

 Höpfner, R., Kutoyants, Y.: On a problem of statistical inference in null recurrent diffusions. Statist. Inference Stoch. Processes 6, 25–42 (2003).

### 7 The Gaussian shift model

**Definition:** A model  $(\Omega, \mathcal{A}, \{P_h : h \in \mathbb{R}^d\})$  is called <u>Gaussian shift experiment  $\mathcal{E}(J)$ </u> if there exists a statistic S,  $\mathbb{R}^d$ -valued, and a deterministic matrix J, symmetric and strictly positive definite, such that for every  $h \in \mathbb{R}^d$ ,

$$\omega \longrightarrow \exp\left(h^{\top}S(\omega) - \frac{1}{2}h^{\top}Jh\right) =: L^{h/0}(\omega)$$

is a version of the likelihood ratio of  $P_h$  with respect to  $P_0$ . Call  $Z := J^{-1}S$  <u>central statistic</u> in  $\mathcal{E}(J)$ .

**Remark:** In a Gaussian shift experiment  $\mathcal{E}(J)$ , all probability measures are equivalent, and we have

$$L^{(h_0+h)/h_0} = L^{(h_0+h)/0}/L^{h_0/0} = \exp\left(h^\top (S-Jh_0) - \frac{1}{2}h^\top Jh\right) \quad \text{for all } h, \, h_0 \text{ in } \mathbb{R}^d \,.$$

Taking expectation  $E_{h_0}(\ldots)$  on both sides we identify (Laplace transform!)

$$\mathcal{L}\left(S - Jh_0 \mid P_{h_0}\right) = \mathcal{N}(0, J)$$

and thus obtain for the central statistic  $Z := J^{-1}S$  and arbitrary  $h_0 \in \mathbb{R}^d$ 

$$\mathcal{L}(Z - h_0 \mid P_{h_0}) = \mathcal{N}(0, J^{-1})$$
 does not depend on  $h_0 \in \mathbb{R}^d$ .

We all estimators  $\kappa$  for the unknown parameter equivariant if laws of estimation errors  $\mathcal{L}(\kappa - h | P_h)$ do not depend on the value of the parameter h. In particular, in a Gaussian shift model  $\mathcal{E}(J)$ , the central statistic Z is an equivariant estimator for the unknown parameter.

Convolution theorem (Boll 1955): In a Gaussian shift experiment  $\mathcal{E}(J)$ , for every equivariant estimator  $\kappa$  for the unknown parameter  $h \in \mathbb{R}^d$ , there is a probability law Q on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  such that

$$\mathcal{L}(\kappa - h \mid P_h) = \mathcal{N}(0, J^{-1}) \star Q \text{ for all } h \in \mathbb{R}^d$$

Convolutions  $\mathcal{N}(0, J^{-1}) \star Q$  are always 'more spread out' than the law  $\mathcal{N}(0, J^{-1})$  itself. Hence in  $\mathcal{E}(J)$ , the central statistic Z is a best concentrated equivariant estimator for the unknown parameter. However, we would like to be able to compare quite arbitrary estimators for the unknown parameter. For this, the main technical (and somewhat difficult) step:

**Lemma:** In a Gaussian shift experiments  $\mathcal{E}(J)$ , for arbitrary estimators  $\eta$  for the unknown parameter  $h \in \mathbb{R}^d$ , there is a sequence  $(Q_n)_n$  of probability measures on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  such that total variation distance between the laws

$$\int_{\{|h| \le n\}} \mathcal{L}(\eta - h \mid P_h) dh \quad \text{and} \quad \mathcal{N}(0, J^{-1}) \star Q_n$$

vanishes as  $n \to \infty$ .

**Definition:** We call a loss function  $\ell : \mathbb{R}^d \to [0, \infty)$  <u>subconvex</u> or bowl-shaped if level sets

$$\left\{ x \in \mathbb{R}^d : \ell(x) \le c \right\} \quad , \quad c \ge 0$$

are convex and symmetric with respect to the origin.

With respect to a loss function  $\ell$ , the <u>risk</u> of an estimator  $\eta$  for the unknown parameter h is the function

$$h \longrightarrow R_{\ell}(\eta, h) := E_h(\ell(\eta - h))$$

which may take the value  $\infty$ .

**Minimax theorem:** In a Gaussian shift experiment  $\mathcal{E}(J)$ , the central statistic Z minimizes the maximal risk with respect to any subconvex loss function  $\ell(\cdot)$ , i.e.: arbitrary estimators  $\eta$  for the unknown parameter  $h \in \mathbb{R}^d$  can be compared to Z via

$$\sup_{h \in \mathbb{R}^d} R_{\ell}(\eta, h) \geq \int_{\mathbb{R}^d} \ell(z) \mathcal{N}(0, J^{-1})(dz) = R_{\ell}(Z, 0) = \sup_{h \in \mathbb{R}^d} R_{\ell}(Z, h) + \sum_{h \in \mathbb{R}^d} R_{\ell}$$

References, proofs, further reading: H. 2014 Section 5.1; mixed normal experiments: Section 6.1.

# 8 Local asymptotic normality (LAN)

Consider a sequence of statistical experiments parametrized by the same parameter set  $\Theta \subset \mathbb{R}^d$  open:

$$\mathcal{E}_n := (\Omega_n, \mathcal{A}_n, \{P_{n,\vartheta} : \vartheta \in \Theta\}) \quad , \quad n \in \mathbb{N} .$$

**Definition:** For  $\vartheta \in \Theta$  fixed, the sequence of experiments  $(\mathcal{E}_n)_n$  is called <u>locally asymptotically</u> normal (LAN) at  $\vartheta$  if there is a sequence of positive real numbers (or matrices in  $\mathbb{R}^{d \times d}$ )

$$\delta_n = \delta_n(\vartheta)$$
 decreasing to 0 as  $n \to \infty$ 

called <u>local scale at  $\vartheta$ </u>, and a matrix in  $\mathbb{R}^{d \times d}$ 

 $J(\vartheta)$  , deterministic, symmetric and strictly positive definite

called <u>limit information at  $\vartheta$ </u>, such that the following holds: for all bounded sequences  $(h_n)_n$  in  $\mathbb{R}^d$ (we always assume that  $\vartheta + \delta_n(\vartheta)h_n$  belongs to  $\Theta$ ), log-likelihood ratios in <u>local models at  $\vartheta$ </u>

$$\mathcal{E}_{n,\vartheta} := \left\{ P_{n\,,\,\vartheta+\delta_n(\vartheta)h} : h \text{ in } \mathbb{R}^d \text{ such that } \vartheta+\delta_n(\vartheta)h \in \Theta \right\} \quad, \quad n \to \infty$$

admit quadratic expansions

$$\log L_n^{(\vartheta+\delta(\vartheta)h_n)/\vartheta} = h_n^{\top} S_n(\vartheta) - \frac{1}{2} h_n^{\top} J_n(\vartheta) h_n + o_{P_{n,\vartheta}}(1)$$

as  $n \to \infty$  where

$$\mathcal{L}\left(S_n(\vartheta), J_n(\vartheta) \mid P_{n,\vartheta}\right) \xrightarrow{w} \mathcal{N}\left(0, J(\vartheta)\right) \otimes \epsilon_{J(\vartheta)} \quad (\text{weakly in } \mathbb{R}^d \times \mathbb{R}^{d \times d} \ ) \ .$$

**Remarks:** 1) Equivalently, we can write the last assertion as

$$\mathcal{L}(S_n(\vartheta) \mid P_{n,\vartheta}) \xrightarrow{w} \mathcal{N}(0, J(\vartheta)) \text{ together with } J_n(\vartheta) = J(\vartheta) + o_{P_{n,\vartheta}}(1)$$

as  $n \to \infty$  (weak convergence to a deterministic object is convergence in probability).

2) Local models at  $\vartheta$  are parametrized by  $h \in ...\mathbb{R}^d...$ : thus, when LAN holds at  $\vartheta$ , the Gaussian shift experiment  $\mathcal{E}(J(\vartheta))$  appears as <u>limit model</u> for  $\mathcal{E}_{n,\vartheta}$  as  $n \to \infty$ . **Definition:** Assuming LAN at  $\vartheta$ , a sequence of estimators for the unknown parameter  $\vartheta \in \Theta$ 

$$T_n$$
 on  $(\Omega_n, \mathcal{A}_n), n \ge 1$ 

is called regular at  $\vartheta$  if there is some probability law  $F = F(\vartheta)$  on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  such that

$$\mathcal{L}\left(\left.\delta_{n}^{-1}(\vartheta)\left(T_{n}-(\vartheta+\delta_{n}(\vartheta)h)\right)\right|P_{n\,,\,\vartheta+\delta_{n}(\vartheta)h}\right) \qquad \stackrel{w}{\longrightarrow} \qquad F$$

(weak convergence in  $\mathbb{R}^d$ , as  $n \to \infty$ ) for every  $h \in \mathbb{R}^d$ .

Thus regular means 'asymptotically equivariant with respect to the local parameter' at  $\vartheta$ . Regularity can be checked by proving joint weak convergence of pairs 'rescaled estimation errors, log-likelihood ratios' under  $P_{n,\vartheta}$  as  $n \to \infty$  to limit laws of suitable structure.

A key tool is 'Le Cam's 3rd lemma', e.g. H. 2014 section 3.1.

**Hájek's convolution theorem:** Assume LAN at  $\vartheta$ , and consider any sequence of estimators  $(T_n)_n$  for the unknown parameter  $\vartheta \in \Theta$  which is regular at  $\vartheta$ .

a) Any limit distribution F arising in the definition above can be written as

$$F = \mathcal{N}(0, J^{-1}(\vartheta)) \star Q$$

for some probability law Q on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ .

b) A sequence of estimators  $(T_n)_n$  is regular and efficient at  $\vartheta$  if and only if

$$\delta_n^{-1}(\vartheta) (T_n - \vartheta) = Z_n(\vartheta) + o_{P_{n,\vartheta}}(1) \text{ as } n \to \infty$$

i.e. rescaled estimation errors under  $\vartheta$  are coupled to the central sequence  $Z_n(\vartheta) = J_n^{-1}(\vartheta)S_n(\vartheta)$  at  $\vartheta$ .

**Local asymptotic minimax theorem:** Assume LAN at  $\vartheta$ , and consider any sequence of estimators  $(T_n)_n$  for the unknown parameter  $\vartheta \in \Theta$  whose rescaled estimation errors at  $\vartheta$  are tight:

$$\mathcal{L}\left(\delta_n^{-1}(\vartheta)\left(T_n-\vartheta\right)\mid P_{n,\vartheta}\right) \ , \ n\geq 1 \ , \quad \text{is tight in } \mathbb{R}^d \text{ as } n\to\infty \ .$$

a) For every loss function  $\ell(\cdot)$  which is continuous, subconvex and bounded:

$$\lim_{c \uparrow \infty} \limsup_{n \to \infty} \sup_{|h| \le c} E_{\vartheta + \delta_n(\vartheta)h} \left( \ell \left( \delta_n^{-1}(\vartheta) \left( \widetilde{\vartheta}_n - (\vartheta + \delta_n(\vartheta)h) \right) \right) \right) \ge \int_{\mathbb{R}^d} \ell(z) \, \mathcal{N}(0, J^{-1}(\vartheta))(dz)$$

where  $J(\vartheta)$  is the limiting information.

b) Sequences  $(T_n)_n$  with the property

$$\delta_n^{-1}(\vartheta) \left(T_n - \vartheta\right) = Z_n(\vartheta) + o_{P_{n,\vartheta}}(1) \text{ as } n \to \infty$$

achieve the local asymptotic minimax bound in a): for every  $0 < c < \infty$ ,

$$\lim_{n \to \infty} \sup_{|h| \le c} E_{\vartheta + \delta_n(\vartheta)h} \left( \ell \left( \delta_n^{-1}(\vartheta) \left( \widetilde{\vartheta}_n - (\vartheta + \delta_n(\vartheta)h) \right) \right) \right) = \int_{\mathbb{R}^d} \ell(z) \, \mathcal{N}(0, J^{-1}(\vartheta))(dz) \, .$$

**Remark:** Starting from any sequence  $(\tilde{T}_n)_n$  of preliminary estimators for the unknown parameter  $\vartheta \in \Theta$  such that rescaled estimation errors at  $\vartheta$  are tight, it is possible (under additional assumptions which often do hold) to <u>construct explicitely</u> a new sequence  $(T_n)_n$  which is efficient in the sense of (the convolution theorem and) the local asymptotic minimax theorem: this is Le Cam's 'one-step correction'.

References, proofs, further reading, extension to LAMN: H. 2014 Sections 7.1, 7.2, 7.3.