# Fisher and Wilks expansions with applications to statistical inference

Vladimir Spokoiny,

WIAS, HU Berlin

# Outline

Let $\boldsymbol{\vartheta}$, a random element $\Theta$,

$\pi(\boldsymbol{\theta})$ a prior density.

The posterior distribution of $\boldsymbol{\vartheta}$ is given by

$$I\!P(A \mid \boldsymbol{Y}) = \frac{\int_A \exp\{L(\boldsymbol{\theta})\}\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}{\int_\Theta \exp\{L(\boldsymbol{\theta})\}\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

Introduce the posterior moments

$$\overline{\boldsymbol{\vartheta}} \stackrel{\text{def}}{=} I\!E(\boldsymbol{\vartheta} \mid \boldsymbol{Y}),$$

$$\mathfrak{S}^2 \stackrel{\text{def}}{=} \text{Cov}(\boldsymbol{\vartheta} \mid \boldsymbol{Y}) \stackrel{\text{def}}{=} I\!E\{(\boldsymbol{\vartheta} - \overline{\boldsymbol{\vartheta}})(\boldsymbol{\vartheta} - \overline{\boldsymbol{\vartheta}})^\top \mid \boldsymbol{Y}\}.$$

## Some references

There is a number of papers in this direction recently appeared:

- [Ghosal et al., 2000, Ghosal and van der Vaart, 2007] for a general theory in the i.i.d. case;
- [Ghosal, 1999], [Ghosal, 2000] for high dimensional linear models;
- [Boucheron and Gassiat, 2009], [Kim, 2006] for some special non-Gaussian models;
- [Shen, 2002], [Bickel and Kleijn, 2012], [Rivoirard and Rousseau, 2012], [Castillo, 2012], [Castillo and Rousseau, 2013] for a semiparametric version of the BvM result for different models;
- [Kleijn and van der Vaart, 2006], [Bunke and Milhaud, 1998], for the misspecified parametric case,
- [Castillo and Rousseau, 2013],
- [Kleijn and van der Vaart, 2012] for a general framework for the BvM result in terms of a stochastic LAN condition

Extensions to nonparametric models with infinite or growing parameter dimension $p$ exist for some special situations:

- [Freedman, 1999] and [Ghosal, 1999, Ghosal, 2000] for linear models
- [Bontemps, 2011] for Gaussian regression,
- [Castillo and Nickl, 2013] for the white noise case;

Below $\pi(\boldsymbol{\theta}) \equiv 1$, an improper non-informative prior.

Yields for any $A \subset \Theta$

$$\mathbb{P}^\circ(A) = \mathbb{P}(A \mid \boldsymbol{Y}) = \frac{\int_A \exp\{L(\boldsymbol{\theta})\} d\boldsymbol{\theta}}{\int_\Theta \exp\{L(\boldsymbol{\theta})\} d\boldsymbol{\theta}} .$$

Quasi-likelihood $\implies$ quasi-posterior.

A general case with a continuous prior density $\pi(\boldsymbol{\theta})$:

$$\boldsymbol{\vartheta} \mid \boldsymbol{Y} \propto \exp\{L(\boldsymbol{\theta})\} \pi(\boldsymbol{\theta}) = \exp\{L_\pi(\boldsymbol{\theta})\}$$

with

$$L_\pi(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) + \log \pi(\boldsymbol{\theta}).$$

So, the case of a general smooth prior can be reduced to the case of a non-informative prior by changing the log-likelihood function.

**Theorem**

*Suppose the conditions of Theorem 19. Let also* $\mathtt{b}(\mathbf{r})$ *from* $(\mathcal{L})$ *satisfies*

$$\mathbf{r}^2\mathtt{b}^2(\mathbf{r}) \ \geq \ \mathtt{x} + 2p + 4z^2(B, \mathbf{x}) + 8\mathbf{r}\,\mathtt{b}(\mathbf{r})\varrho(\mathbf{r}, \mathbf{x}), \qquad \mathbf{r} \geq \mathbf{r}_0, \tag{1}$$

*with* $\varrho(\mathbf{r}, \mathbf{x})$ *from* (14). *Then it holds on a random set* $\Omega(\mathbf{x})$ *of probability at least* $1 - 5\mathrm{e}^{-\mathtt{x}}$

$$I\!P\big(\boldsymbol{\vartheta} \notin \Theta_0(\mathbf{r}_0) \,\big|\, \boldsymbol{Y}\big) \ \leq \ \mathrm{e}^{-\mathtt{x}}.$$

The bound (1) is very similar to the bound for the MLE concentration. It can be spelled out as the condition that

▶ $\mathbf{r}_0^2 \geq 2p + \mathtt{x} + 4\mathfrak{z}^2(B, \mathbf{x})$,

▶ $\mathtt{b}(\mathbf{r}_0) \approx 1$, and

▶ $\mathbf{r}\mathtt{b}(\mathbf{r})$ grows with $\mathbf{r}$.

Define

$$\breve{\boldsymbol{\theta}} = \boldsymbol{\theta}^* + D_0^{-2}\nabla L(\boldsymbol{\theta}^*) = \boldsymbol{\theta}^* + D_0^{-1}\boldsymbol{\xi}.$$

The Fisher result implies

$$\|D_0(\widetilde{\boldsymbol{\theta}} - \breve{\boldsymbol{\theta}})\| \leq \diamondsuit(\mathtt{r}_0, \mathtt{x}).$$

**Theorem**

*On* $\Omega(\mathtt{x})$

$$\|D_0(\overline{\boldsymbol{\vartheta}} - \breve{\boldsymbol{\theta}})\|^2 \leq 4\Delta(\mathtt{r}_0, \mathtt{x}) + 4\mathrm{e}^{-\mathtt{x}},$$

$$\left\|I_p - D_0\mathfrak{S}^2 D_0\right\|_\infty \leq 4\Delta(\mathtt{r}_0, \mathtt{x}) + 4\mathrm{e}^{-\mathtt{x}}.$$

$$\breve{\boldsymbol{\theta}} = \boldsymbol{\theta}^* + D_0^{-2}\nabla L(\boldsymbol{\theta}^*) = \boldsymbol{\theta}^* + D_0^{-1}\boldsymbol{\xi}.$$

**Theorem**

*For any* $\boldsymbol{\lambda} \in \mathbb{R}^p$ *with* $\|\boldsymbol{\lambda}\|^2 \leq p$

$$\left|\log \mathbb{E}\Big[\exp\big\{\boldsymbol{\lambda}^\top D_0(\boldsymbol{\vartheta} - \breve{\boldsymbol{\theta}})\big\} \,\big|\, \boldsymbol{Y}\Big] - \|\boldsymbol{\lambda}\|^2/2\right| \leq 2\Delta(\mathbf{r}_0, \mathbf{x}) + 3\mathrm{e}^{-\mathbf{x}},$$

*and for any measurable set* $A \subset \mathbb{R}^p$

$$\mathbb{P}\big(D_0(\boldsymbol{\vartheta} - \breve{\boldsymbol{\theta}}) \in A \,\big|\, \boldsymbol{Y}\big) \geq \exp\big\{-2\Delta(\mathbf{r}_0, \mathbf{x}) - 3\mathrm{e}^{-\mathbf{x}}\big\}\mathbb{P}\big(\boldsymbol{\gamma} \in A\big) - \mathrm{e}^{-\mathbf{x}},$$

$$\mathbb{P}\big(D_0(\boldsymbol{\vartheta} - \breve{\boldsymbol{\theta}}) \in A \,\big|\, \boldsymbol{Y}\big) \leq \exp\big\{2\Delta(\mathbf{r}_0, \mathbf{x}) + 2\mathrm{e}^{-\mathbf{x}}\big\}\mathbb{P}\big(\boldsymbol{\gamma} \in A\big) + \mathrm{e}^{-\mathbf{x}}.$$

► All statements of Theorem 1 require "$\Delta(r_0, x)$ is small".

► The BvM result is stated under essentially the same list of conditions as the frequentist results of Fisher and Wilks Theorems.

► The normal approximation of the posterior is entirely based on the smoothness properties of the likelihood function

► No any asymptotic arguments like weak convergence or convergence in probability, or the Central Limit Theorem.

► The results continue to hold if $\breve{\boldsymbol{\theta}}$ is replaced by any efficient estimate $\widehat{\boldsymbol{\theta}}$, e.g. by the MLE $\widetilde{\boldsymbol{\theta}}$, satisfying $\|D_0(\widehat{\boldsymbol{\theta}} - \breve{\boldsymbol{\theta}})\| \leq r_0$ with a dominating probability.

**Steps: Local Gaussian approximation of the posterior**

Remind $D_0^2 = -\nabla^2 I\!EL(\boldsymbol{\theta}^*)$, $\boldsymbol{\xi} = D_0^{-1} \nabla L(\boldsymbol{\theta}^*)$, and

$$\breve{\boldsymbol{\theta}} = \boldsymbol{\theta}^* + D_0^{-1}\boldsymbol{\xi} = \boldsymbol{\theta}^* + D_0^{-2}\nabla L(\boldsymbol{\theta}^*)$$

Local approximation: on $\Omega(\mathbf{x})$, for $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \frac{1}{2}\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2$

$$\left| L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \right| \leq \Delta(\mathbf{r}_0, \mathbf{x}), \quad \boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0). \tag{2}$$

For any nonnegative function $f$, it holds

$$\int_{\Theta_0(\mathbf{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big)\, d\boldsymbol{\theta}$$

$$\leq \mathrm{e}^{\Delta(\mathbf{r}_0, \mathbf{x})} \int_{\Theta_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big)\, d\boldsymbol{\theta}.$$

$$\int_{\Theta_0(\mathbf{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big)\, d\boldsymbol{\theta}$$

$$\geq \mathrm{e}^{-\Delta(\mathbf{r}_0, \mathbf{x})} \int_{\Theta_0(\mathbf{r}_0)} \exp\{\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big)\, d\boldsymbol{\theta}.$$

The main benefit: $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is quadratic in $\boldsymbol{\theta}$ and thus

$$\exp \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \exp\{\boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2\}$$

is proportional to the density of a Gaussian distribution.

More precisely, define

$$m(\boldsymbol{\xi}) \stackrel{\text{def}}{=} -\|\boldsymbol{\xi}\|^2/2 + \log(\det D_0) - p \log(\sqrt{2\pi}).$$

Then

$$m(\boldsymbol{\xi}) + \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \;=\; -\|D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\|^2/2 + \log(\det D_0) - p \log(\sqrt{2\pi}) \qquad (3)$$

is (conditionally on $\boldsymbol{Y}$) the log-density of the normal law $\mathcal{N}(\breve{\boldsymbol{\theta}}, D_0^{-2})$ with the mean $\breve{\boldsymbol{\theta}} = D_0^{-1}\boldsymbol{\xi} + \boldsymbol{\theta}^*$ and the covariance matrix $D_0^{-2}$.

**Theorem**

*For any nonnegative function $f(\cdot)$ on $\mathbb{R}^p$, it holds on $\Omega(\mathbf{r}_0, \mathbf{x})$*

$$\mathbb{E}^\circ\left[f\left(D_0(\boldsymbol{\vartheta} - \check{\boldsymbol{\theta}})\right)\mathbb{1}_{\mathbf{r}_0}\right] \leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\}\,\mathbb{E}f(\boldsymbol{\gamma}), \tag{4}$$

*where*

$$\Delta^+(\mathbf{r}_0, \mathbf{x}) = 2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0),$$

$$\nu(\mathbf{r}_0) \stackrel{\text{def}}{=} -\log \mathbb{P}^\circ\left(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0\right).$$

*If $\mathbf{r}_0^2 \geq z^2(B, \mathbf{x}) + p + 2\mathbf{x}$, then on $\Omega(B, \mathbf{x})$, it holds*

$$\nu(\mathbf{r}_0) \leq 2\mathrm{e}^{-\mathbf{x}}$$

$$\Delta^+(\mathbf{r}_0, \mathbf{x}) \leq 2\Delta(\mathbf{r}_0, \mathbf{x}) + 2\mathrm{e}^{-\mathbf{x}}.$$

We use that $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2/2$ is proportional to the density of a Gaussian distribution. More precisely, define

$$m(\boldsymbol{\xi}) \stackrel{\text{def}}{=} -\|\boldsymbol{\xi}\|^2/2 + \log(\det D_0) - p\log(\sqrt{2\pi}).$$

Then

$$m(\boldsymbol{\xi}) + \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = -\|D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\|^2/2 + \log(\det D_0) - p\log(\sqrt{2\pi}) \tag{5}$$

is (conditionally on $\boldsymbol{Y}$) the log-density of the normal law with the mean $\breve{\boldsymbol{\theta}} = D_0^{-1}\boldsymbol{\xi} + \boldsymbol{\theta}^*$ and the covariance matrix $D_0^{-2}$. Change of variables $\boldsymbol{u} = D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})$ implies by (5) for any nonnegative function $f$ that

$$\int_{\Theta_0(\mathfrak{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + m(\boldsymbol{\xi})\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big)\, d\boldsymbol{\theta}$$

$$\leq e^{\Delta(\mathfrak{r}_0, \mathbf{x})} \exp\{m(\boldsymbol{\xi})\} \int \exp\{\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big)\, d\boldsymbol{\theta}$$

$$= e^{\Delta(\mathfrak{r}_0, \mathbf{x})} \int \phi(\boldsymbol{u})\, f(\boldsymbol{u})\, d\boldsymbol{u} = e^{\Delta(\mathfrak{r}_0, \mathbf{x})}\, \mathbb{E}f(\boldsymbol{\gamma}). \tag{6}$$

Similarly, for any nonnegative function $f$, it follows by change of variables $\boldsymbol{u} = D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})$ and $D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) = \boldsymbol{u} + \boldsymbol{\xi}$ that

$$\exp\{m(\boldsymbol{\xi})\} \int \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f\big(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\big) \, \mathbb{1}\{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \le \mathtt{r}_0\} d\boldsymbol{\theta}$$

$$\ge \exp\{-\Delta(\mathtt{r}_0, \mathbf{x})\} \int \phi(\boldsymbol{u}) f(\boldsymbol{u}) \, \mathbb{1}\{\|\boldsymbol{u} + \boldsymbol{\xi}\| \le \mathtt{r}_0\} d\boldsymbol{u}. \tag{7}$$

A special case of (7) with $f(\boldsymbol{u}) \equiv 1$ implies by definition of $\nu(\mathtt{r}_0)$:

$$\exp\{m(\boldsymbol{\xi})\} \int_{\Theta_0(\mathtt{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} \, d\boldsymbol{\theta} \ge \exp\{-\Delta(\mathtt{r}_0, \mathbf{x}) - \nu(\mathtt{r}_0)\}. \tag{8}$$

Now we are prepared to finalize the proof. (6) and (8) imply on $\Omega(\mathbf{r}_0, \mathbf{x})$

$$\frac{\int_{\Theta_0(\mathbf{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f(D_0(\boldsymbol{\theta} - \boldsymbol{\breve{\theta}}))\, d\boldsymbol{\theta}}{\int \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}\, d\boldsymbol{\theta}} \; \leq \; \exp\{2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0)\} \; I\!\!E f(\boldsymbol{\gamma})$$

and (4) follows. As $\|\boldsymbol{\xi}\| \leq z(B, \mathbf{x})$ on $\Omega(B, \mathbf{x})$ and $\mathbf{r}_0 \geq z(B, \mathbf{x}) + z(p, \mathbf{x})$,

$$\nu(\mathbf{r}_0) \; = \; -\log I\!\!P^{\circ}\big(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0\big) \leq -\log I\!\!P\big(\|\boldsymbol{\gamma}\| \leq z(p, \mathbf{x})\big) \leq 2\mathrm{e}^{-\mathbf{x}},$$

---

**Lemma**

*For each* $\mathbf{x}$ *and for* $\boldsymbol{\gamma} \sim \mathcal{N}(0, I_p)$

$$I\!\!P\big(\|\boldsymbol{\gamma}\| \geq z(p, \mathbf{x})\big) \; \leq \; \exp(-\mathbf{x}), \qquad I\!\!P\big(\|\boldsymbol{\gamma}\| \leq z_1(p, \mathbf{x})\big) \; \leq \; \exp(-\mathbf{x}),$$

*where*

$$z^2(p, \mathbf{x}) \stackrel{\text{def}}{=} p + \sqrt{6.6p\mathbf{x}} \vee (6.6\mathbf{x}), \qquad z_1^2(p, \mathbf{x}) \stackrel{\text{def}}{=} p - 2\sqrt{p\mathbf{x}}.$$

The next important step in our analysis is to check that $\boldsymbol{\vartheta}$ concentrates in a small vicinity $\Theta_0 = \Theta_0(\mathbf{r}_0)$ of the central point $\boldsymbol{\theta}^*$ with a properly selected $\mathbf{r}_0$. The concentration properties of the posterior will be described by using the random quantity

$$\rho(\mathbf{r}_0) \stackrel{\text{def}}{=} \frac{\int_{\Theta \setminus \Theta_0} \exp\{L(\boldsymbol{\theta})\} d\boldsymbol{\theta}}{\int_{\Theta_0} \exp\{L(\boldsymbol{\theta})\} d\boldsymbol{\theta}} .$$

Obviously $I\!\!P\{\boldsymbol{\vartheta} \notin \Theta_0(\mathbf{r}_0) \,\big|\, \boldsymbol{Y}\} \leq \rho(\mathbf{r}_0)$. Therefore, small values of $\rho(\mathbf{r}_0)$ indicate a small posterior probability of the set $\Theta \setminus \Theta_0$. The proof only uses condition $(\mathcal{L})$ and the fact that there exists a random set $\Omega(\mathbf{x})$ of probability at least $1 - e^{-\mathbf{x}}$ such that

$$\big|\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\big| \leq \mathbf{r}\, \varrho(\mathbf{r}, \mathbf{x}) \tag{9}$$

for $\mathbf{r} = \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$; cf. the proof of Theorem 19.

Let $\mathbf{b}_0 = \mathbf{b}(\mathbf{r}_0)$ and for the sequence $\mathbf{b}_k = 2^{-k}\mathbf{b}_0$, the radii $\mathbf{r}_0 < \mathbf{r}_1 < \ldots$ be defined by the condition $\mathbf{b}(\mathbf{r}) \geq \mathbf{b}_k > 0$ for $\mathbf{r}_k \leq \mathbf{r} < \mathbf{r}_{k+1}$ for all $k \geq 0$ with $\mathbf{b}(\mathbf{r})$ from $(\mathcal{L})$.

**Theorem**

*Suppose the conditions* $(\mathcal{L})$, $(ED_0)$, *and* $(ED_2)$. *If* $\mathbf{b}(\mathbf{r})$ *from* $(\mathcal{L})$ *satisfies*

$$\mathbf{r}^2 \mathbf{b}^2(\mathbf{r}) \geq \mathbf{x} + 2p + 4z^2(B, \mathbf{x}) + 8\mathbf{r}\,\mathbf{b}(\mathbf{r})\varrho(\mathbf{r}, \mathbf{x}), \qquad \mathbf{r} \geq \mathbf{r}_0, \qquad (10)$$

*then it holds on a set* $\Omega(\mathbf{x})$ *of probability at least* $1 - 4\mathrm{e}^{-\mathbf{x}}$

$$\rho(\mathbf{r}_0) \stackrel{\text{def}}{=} \frac{\int_{\Theta \setminus \Theta_0} \exp\{L(\boldsymbol{\theta})\}d\boldsymbol{\theta}}{\int_{\Theta_0} \exp\{L(\boldsymbol{\theta})\}d\boldsymbol{\theta}} \leq 2\exp\{-\mathbf{x} + \Delta^+(\mathbf{r}_0, \mathbf{x})\} \qquad (11)$$

*with* $\Delta^+(\mathbf{r}_0, \mathbf{x}) \leq 2\Delta(\mathbf{r}_0, \mathbf{x}) + 2\mathrm{e}^{-\mathbf{x}}$.

Suppose that $\mathbf{b}_0 = \mathbf{b}(\mathbf{r}_0)$ is close to one and $\varrho(\mathbf{r}, \mathbf{x})$ small. Condition (10) requires that

$$\mathbf{r}_0^2 > 4z^2(B, \mathbf{x}) + 2p + \mathbf{x}$$

and the value $\mathbf{r}\mathbf{b}(\mathbf{r})$ grows with $\mathbf{r}$.

Use the decomposition

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = I\!\!E L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) + \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*).$$

$$= I\!\!E L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \boldsymbol{\xi}^\top D_0 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D_0 (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top.$$

Condition $(\mathcal{L})$ for the expected negative log-likelihood implies

$$-I\!\!E L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \left| D_0 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right|^2 \mathtt{b}_k / 2$$

for each $k \geq 0$ and any $\boldsymbol{\theta} \in \Theta_0(\mathtt{r}_{k+1}) \setminus \Theta_0(\mathtt{r}_k)$. The bound (9) implies on $\Omega(\mathtt{x})$

$$\left| \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D_0 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right| \leq \mathtt{r}_{k+1} \, \varrho(\mathtt{r}_{k+1}, \mathtt{x}), \qquad \boldsymbol{\theta} \in \Theta_0(\mathtt{r}_{k+1}) \setminus \Theta_0(\mathtt{r}_k),$$

Represent

$$\rho(\mathtt{r}_0) \stackrel{\text{def}}{=} \frac{\int_{\Theta \setminus \Theta_0} \exp\{L(\boldsymbol{\theta})\} d\boldsymbol{\theta}}{\int_{\Theta_0} \exp\{L(\boldsymbol{\theta})\} d\boldsymbol{\theta}} = \frac{\exp\{m(\boldsymbol{\xi})\} \int_{\Theta \setminus \Theta_0} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} d\boldsymbol{\theta}}{\exp\{m(\boldsymbol{\xi})\} \int_{\Theta_0} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} d\boldsymbol{\theta}}.$$

By the change of variables $\boldsymbol{\gamma} = D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$, it follows for each $k$

$$\exp\{m(\boldsymbol{\xi})\} \int_{\Theta_0(\mathbf{r}_{k+1}) \setminus \Theta_0(\mathbf{r}_k)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} d\boldsymbol{\theta}$$

$$\leq \exp\{\mathbf{r}_{k+1} \, \varrho(\mathbf{r}_{k+1}, \mathbf{x}) - \|\boldsymbol{\xi}\|^2 / 2\} \frac{1}{(2\pi)^{p/2}} \int_{\|\boldsymbol{\gamma}\| \geq \mathbf{r}_k} \exp\left\{-\frac{\mathbf{b}_k \|\boldsymbol{\gamma}\|^2}{2} + \boldsymbol{\xi}^\top \boldsymbol{\gamma}\right\} d\boldsymbol{\gamma}.$$

Next,

$$\frac{1}{(2\pi)^{p/2}} \int_{\|\boldsymbol{\gamma}\| \geq \mathbf{r}_k} \exp\left(-\frac{\mathbf{b}_k \|\boldsymbol{\gamma}\|^2}{2} + \boldsymbol{\xi}^\top \boldsymbol{\gamma}\right) d\boldsymbol{\gamma}$$

$$\leq \mathbf{b}_k^{-p/2} \exp\left(\frac{\|\boldsymbol{\xi}\|^2}{2\mathbf{b}_k}\right) I\!\!P^\circ\left(\|\boldsymbol{\gamma} + \mathbf{b}_k^{-1/2} \boldsymbol{\xi}\| \geq \mathbf{b}_k^{1/2} \mathbf{r}_k\right)$$

$$\leq \mathbf{b}_k^{-p/2} \exp\left(\frac{\|\boldsymbol{\xi}\|^2}{\mathbf{b}_k} - \frac{1}{4} \mathbf{b}_k \mathbf{r}_k^2 + \frac{p}{2}\right). \tag{12}$$

Here we have used the bound for a standard normal vector $\boldsymbol{\gamma}$ and $\boldsymbol{u} = \mathbf{b}_k^{-1/2} \boldsymbol{\xi} \in I\!\!R^p$. (8) and (12) imply (11).

Now the bound $\|\boldsymbol{\xi}\| \leq z(B, \mathtt{x})$ holding with a dominating probability and (10) imply

$$\sum_{k=0}^{\infty} \exp\{m(\boldsymbol{\xi})\} \int_{\Theta_0(\mathtt{r}_{k+1}) \setminus \Theta_0(\mathtt{r}_k)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} d\boldsymbol{\theta}$$

$$\leq \sum_{k=0}^{\infty} \exp\left(\frac{\|\boldsymbol{\xi}\|^2}{\mathtt{b}_k} - \frac{1}{4}\mathtt{b}_k \mathtt{r}_k^2 + \frac{p}{2} \log(e/\mathtt{b}_k) + \mathtt{r}_{k+1}\, \varrho(\mathtt{r}_{k+1}, \mathtt{x})\right)$$

$$\leq \sum_{k=0}^{\infty} \exp(-\mathtt{x}/\mathtt{b}_k) \leq 2e^{-\mathtt{x}}$$

and (11) follows in view of $\mathtt{b} \log(e/\mathtt{b}) \leq 1$ for $\mathtt{b} \leq 1$.

**Theorem**

*Suppose* (2) *for* $\mathbf{r} = \mathbf{r}_0$ *and* (11). *Then for any nonnegative function* $f(\cdot)$ *on* $\mathbb{R}^p$ , *it holds on* $\Omega(\mathbf{x})$

$$\mathbb{E}^{\circ}\left\{ f\left(D_0(\boldsymbol{\vartheta} - \boldsymbol{\check{\theta}})\right) \mathbb{1}_{\mathbf{r}_0} \right\} \geq \exp\left\{-\Delta^-(\mathbf{r}_0, \mathbf{x})\right\} \mathbb{E}\left\{ f(\boldsymbol{\gamma}) \, \mathbb{1}\left(\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0\right) \right\},$$

*where*

$$\Delta^-(\mathbf{r}_0, \mathbf{x}) = \Delta^+(\mathbf{r}_0, \mathbf{x}) + \rho(\mathbf{r}_0).$$

On the set $\Omega(\mathbf{x})$, it holds by (6) with $f(\cdot) = 1$:

$$
\begin{aligned}
\int \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}\, d\boldsymbol{\theta} &\leq \int_{\Theta_0} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}\, d\boldsymbol{\theta} + \int_{\Theta \setminus \Theta_0} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}\, d\boldsymbol{\theta} \\
&\leq \{1 + \rho(\mathbf{r}_0)\} \int_{\Theta_0} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}\, d\boldsymbol{\theta} \\
&\leq \{1 + \rho(\mathbf{r}_0)\} \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) + \nu(\mathbf{r}_0)\} \\
&\leq \exp\{\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) + \nu(\mathbf{r}_0) + \rho(\mathbf{r}_0)\}.
\end{aligned}
$$

This and the bound (7) imply

$$
\frac{\exp\{m(\boldsymbol{\xi})\} \int_{\Theta_0(\mathbf{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} f(D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}}))\, d\boldsymbol{\theta}}{\exp\{m(\boldsymbol{\xi})\} \int \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\}\, d\boldsymbol{\theta}}
$$

$$
\geq \frac{\exp\{-\Delta(\mathbf{r}_0, \mathbf{x})\} \int \phi(\boldsymbol{u}) f(\boldsymbol{u})\, \mathbb{1}\{\|\boldsymbol{u} + \boldsymbol{\xi}\| \leq \mathbf{r}_0\}\, d\boldsymbol{u}}{\exp\{\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0) + \rho(\mathbf{r}_0)\}}
$$

$$
\geq \exp\{-\Delta^-(\mathbf{r}_0, \mathbf{x})\}\, \mathbb{E}\big[f(\boldsymbol{\gamma})\, \mathbb{1}\{\|\boldsymbol{\gamma} + \boldsymbol{\xi}\| \leq \mathbf{r}_0\}\big].
$$

Define $\mathcal{C}^\circ(A) = \big\{ \boldsymbol{\theta} \colon D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}}) \in A \big\}$. Then

$$\mathbb{P}\big(\mathcal{C}^\circ(A) \,\big|\, \boldsymbol{Y}\big) \approx \mathbb{P}(\boldsymbol{\gamma} \in A) \pm \mathtt{C}\,\Delta(\mathbf{r}_0, \mathbf{x}).$$

Unfortunately, the quantities $\breve{\boldsymbol{\theta}}$ and $D_0^2$ are unknown and cannot be used for building the elliptic credible sets.

A natural question: empirical counterparts.

**Theorem**

*Let a vector $\widehat{\boldsymbol{\theta}}$ and a symmetric matrix $\widehat{D}$ fulfill*

$$\|D_0(\widehat{\boldsymbol{\theta}} - \breve{\boldsymbol{\theta}})\| \,\leq\, \beta, \qquad \widehat{D}^2 \,\leq\, a^2 D_0^2, \qquad \operatorname{tr}\big(D_0^{-1}\widehat{D}^2 D_0^{-1} - \boldsymbol{I}_p\big)^2 \,\leq\, \epsilon^2.$$

*Then with $\tau = \frac{1}{2}\sqrt{a^2\beta^2 + \epsilon^2}$, it holds on a random set $\Omega(\mathbf{x})$ of probability $1 - 5\mathrm{e}^{-\mathbf{x}}$*

$$\mathbb{P}\big(\widehat{D}(\boldsymbol{\vartheta} - \widehat{\boldsymbol{\theta}}) \in A \,\big|\, \boldsymbol{Y}\big) \,\geq\, \exp\big(-2\Delta(\mathbf{r}_0, \mathbf{x}) - 3\mathrm{e}^{-\mathbf{x}}\big)\big\{ \mathbb{P}\big(\boldsymbol{\gamma} \in A\big) - \tau \big\} - \mathrm{e}^{-\mathbf{x}},$$

$$\mathbb{P}\big(\widehat{D}(\boldsymbol{\vartheta} - \widehat{\boldsymbol{\theta}}) \in A \,\big|\, \boldsymbol{Y}\big) \,\leq\, \exp\big(2\Delta(\mathbf{r}_0, \mathbf{x}) + 2\mathrm{e}^{-\mathbf{x}}\big)\big\{ \mathbb{P}\big(\boldsymbol{\gamma} \in A\big) + \tau \big\} + \mathrm{e}^{-\mathbf{x}}.$$

Denote $U = \widehat{D} D_0^{-1}$ and $\boldsymbol{\eta} = D_0(\boldsymbol{\vartheta} - \breve{\boldsymbol{\theta}})$, and $\boldsymbol{\beta} = D_0(\widehat{\boldsymbol{\theta}} - \breve{\boldsymbol{\theta}})$. Then

$$IP\big(\widehat{D}(\boldsymbol{\vartheta} - \widehat{\boldsymbol{\theta}}) \in A \,\big|\, \boldsymbol{Y}\big) = IP\big(U(\boldsymbol{\eta} - \boldsymbol{\beta}) \in A \,\big|\, \boldsymbol{Y}\big) \approx IP\big(U(\boldsymbol{\gamma} - \boldsymbol{\beta}) \in A \,\big|\, \boldsymbol{Y}\big).$$

Now the result follows from Theorem 1 and

**Lemma**

*Let* $IP_0 = \mathcal{N}(0, \boldsymbol{I}_p)$ *and* $IP_1 = \mathcal{N}(\boldsymbol{\beta}, (U^\top U)^{-1})$ *some non-degenerated matrix* $U$. *If*

$$\|U^\top U - \boldsymbol{I}_p\|_\infty \leq \boldsymbol{\epsilon} \leq 1/2,$$

*then* $\mathcal{K}(IP_0, IP_1) = -I\!\!E_0 \log \frac{dIP_1}{dIP_0}$ *fulfills*

$$2\mathcal{K}(IP_0, IP_1) \leq \operatorname{tr}(U^\top U - \boldsymbol{I}_p)^2 + (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2 \leq \boldsymbol{\epsilon}^2\, p + (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2.$$

*For any measurable set* $A \subset I\!\!R^p$, *it holds with* $\boldsymbol{\gamma} \sim \mathcal{N}(0, \boldsymbol{I}_p)$

$$\big|IP_0(A) - IP_1(A)\big| = \big|IP\big(\boldsymbol{\gamma} \in A\big) - IP\big(U(\boldsymbol{\gamma} - \boldsymbol{\beta}) \in A\big)\big| \leq \sqrt{\mathcal{K}(IP_0, IP_1)/2}.$$

**Proof**

It holds

$$2 \log \frac{d\mathbb{P}_1}{d\mathbb{P}_0}(\boldsymbol{\gamma}) = \log \det(U^\top U) - (\boldsymbol{\gamma} - \boldsymbol{\beta})^\top U^\top U (\boldsymbol{\gamma} - \boldsymbol{\beta}) + \|\boldsymbol{\gamma}\|^2$$

with $\boldsymbol{\gamma}$ standard normal and

$$2\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) = -2\mathbb{E}_0 \log \frac{d\mathbb{P}_1}{d\mathbb{P}_0} = -\log \det(U^\top U) + \operatorname{tr}(U^\top U - \boldsymbol{I}_p) + \boldsymbol{\beta}^\top U^\top U \boldsymbol{\beta}.$$

Let $a_j$ be the $j$ th eigenvalue of $U^\top U - \boldsymbol{I}_p$. $\|U^\top U - \boldsymbol{I}_p\|_\infty \leq \boldsymbol{\epsilon} \leq 1/2$ yields $|a_j| \leq 1/2$ and

$$2\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) = \boldsymbol{\beta}^\top U^\top U \boldsymbol{\beta} + \sum_{j=1}^p \{a_j - \log(1 + a_j)\} \leq (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2 + \sum_{j=1}^p a_j^2$$

$$\leq (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2 + \operatorname{tr}(U^\top U - \boldsymbol{I}_p)^2 \leq (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2 + \boldsymbol{\epsilon}^2 p.$$

This implies by Pinsker's inequality

$$\sup_A |\mathbb{P}_0(A) - \mathbb{P}_1(A)| \leq \sqrt{\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1)/2}.$$

Define

$$\mathcal{C}(A_\alpha) = \big\{ \boldsymbol{\theta} \colon \widehat{D}(\boldsymbol{\theta} - \widehat{\boldsymbol{\theta}}) \in A_\alpha \big\},$$

where $\widehat{D}^2 \approx D_0^2$ and $\widehat{\boldsymbol{\theta}} \approx \breve{\boldsymbol{\theta}} = \boldsymbol{\theta}^* + D_0^{-1}\boldsymbol{\xi} \approx \widetilde{\boldsymbol{\theta}}$. Then

$$\mathbb{P}^\circ\big\{ \mathcal{C}(A_\alpha) \big\} \approx \mathbb{P}(\boldsymbol{\gamma} \in A_\alpha) \pm \mathtt{C}\,\Delta(\mathbf{r}_0, \mathbf{x}).$$

$\mathcal{C}(A_\alpha)$ is completely data-based, can be constructed by Bayesian simulations and $\mathbb{P}^\circ\big\{ \mathcal{C}(A_\alpha) \big\} \approx \alpha$!

Question: can one use $\mathcal{C}(A_\alpha)$ as a frequentist confidence set?

The construction of $\mathcal{C}(A_\alpha)$ perfectly matches the usual frequentist asymptotic CS.

- Under PA $\mathcal{C}(A_\alpha)$ is an asymptotic $\alpha$-CS.

- If PA-PW, the CS $\mathcal{C}(A_\alpha)$ can be totally wrong, cf. [Cox, 1993] or [Kleijn and van der Vaart, 2012].

# Outline

Data $\boldsymbol{Y}$ with DGP $\boldsymbol{Y} \sim \mathbb{P}$.

SPA: $\mathbb{P} \in (\mathbb{P}_{\boldsymbol{\theta}, \boldsymbol{\eta}})$, probably misspecified.

$\boldsymbol{\theta}$, target, $\dim(\boldsymbol{\theta}) = p$, $\quad \boldsymbol{\eta}$, nuisance, $\dim(\boldsymbol{\eta}) = q$, $\quad p^* = p + q$.

Goal: inference on $\boldsymbol{\theta}$.

Examples in mind:

- an inverse problem with error in operator;
  $\boldsymbol{Y} = A\boldsymbol{\theta} + \boldsymbol{\varepsilon}$, observed $\boldsymbol{Y}$ and $\widehat{A}$, operator $A$ as nuisance;

- transformation models $\Lambda Y = f(X) + \epsilon$: the transfer $\Lambda$ or regression function $\boldsymbol{f}$ as nuisance;

- Hidden Markov Chains $Y_t \sim P_{f(X_t, \boldsymbol{\theta})}$: the whole hidden path $X_t$ as nuisance.

- Error-in-variable regression $Y_i = f(X_i) + \epsilon_i$, $Z_i = X_i + \xi_i$: the whole unobserved design $\boldsymbol{X}$ as nuisance.

SPA : $\qquad \boldsymbol{Y} \sim \boldsymbol{P} \in \big( \boldsymbol{P}_{\boldsymbol{\theta}, \boldsymbol{\eta}}, \; \boldsymbol{\theta} \in \Theta, \boldsymbol{\eta} \in H \big)$

Log-likelihood: $\quad L(\boldsymbol{\theta}, \boldsymbol{\eta}) = \dfrac{d\boldsymbol{P}_{\boldsymbol{\theta}, \boldsymbol{\eta}}}{d\boldsymbol{\mu}_0}(\boldsymbol{Y})$

Profile MLE: $\quad \widetilde{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \, \underset{\boldsymbol{\eta}}{\max} \, L(\boldsymbol{\theta}, \boldsymbol{\eta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \, \breve{L}(\boldsymbol{\theta}), \qquad \breve{L}(\boldsymbol{\theta}) = \underset{\boldsymbol{\eta}}{\max} \, L(\boldsymbol{\theta}, \boldsymbol{\eta}).$

Murphy, van der Vaart (2000), Kosorok (2005, 2008): Under PA $\boldsymbol{P} = \boldsymbol{P}_{\boldsymbol{\theta}^*, \boldsymbol{\eta}^*}$, the pMLE $\widetilde{\boldsymbol{\theta}}$ is

- ■ root- $n$ consistent and normal
- ■ semiparametrically efficient
- ■ $2\breve{L}(\widetilde{\boldsymbol{\theta}}) - 2\breve{L}(\boldsymbol{\theta}^*) \xrightarrow{w} \chi_p^2$, where $p = \dim(\Theta)$.

Limitations:

- ■ hard optimization problem, often unfeasible
- ■ SPA is crucial but questionable
- ■ large sample asymptotics

$(\boldsymbol{\theta}, \boldsymbol{\eta})$ -setup:

$$\boldsymbol{Y} = \Psi^\top \boldsymbol{\theta}^* + \Phi^\top \boldsymbol{\eta}^* + \boldsymbol{\varepsilon},$$

where $\Psi$ is $p \times n$ matrix of essential factors $\boldsymbol{\psi}_1, \ldots, \boldsymbol{\psi}_p$, $\Phi$ is $q \times n$ -matrix of nuisance factors $\boldsymbol{\phi}_1, \ldots, \boldsymbol{\phi}_q$.

$\boldsymbol{\upsilon}$ -setup:

$$\boldsymbol{Y} = \Upsilon^\top \boldsymbol{\upsilon}^* + \boldsymbol{\varepsilon}$$

with $p^*$ factors $(\boldsymbol{\psi}_j), (\boldsymbol{\phi}_m)$, and the target of estimation is a linear mapping $\boldsymbol{\theta}^* = P\boldsymbol{\upsilon}^*$ for a given projector $P : I\!R^{p^*} \to I\!R^p$.

$\boldsymbol{v}$ -setup:

$$\boldsymbol{Y} = \Upsilon^{\top}\boldsymbol{v}^* + \boldsymbol{\varepsilon} = \Psi^{\top}\boldsymbol{\theta}^* + \Phi^{\top}\boldsymbol{\eta}^* + \boldsymbol{\varepsilon}, \qquad I\!\!E\boldsymbol{\varepsilon} = 0, \mathrm{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n.$$

Target: $\boldsymbol{\theta}^* = P\boldsymbol{v}^*$.

Profile qMLE 1: $\qquad \widetilde{\boldsymbol{\theta}} = P\widetilde{\boldsymbol{v}} = P(\Upsilon\Upsilon^{\top})^{-1}\Upsilon\boldsymbol{Y} = S\boldsymbol{Y}, \qquad\qquad S = P(\Upsilon\Upsilon^{\top})^{-1}\Upsilon.$

Profile qMLE 2: $\qquad \widetilde{\boldsymbol{\theta}} \stackrel{\mathrm{def}}{=} \mathop{\mathrm{argmax}}\limits_{\boldsymbol{\theta}} \breve{L}(\boldsymbol{\theta}), \qquad\qquad \breve{L}(\boldsymbol{\theta}) \stackrel{\mathrm{def}}{=} \sup\limits_{\boldsymbol{v}:\ P\boldsymbol{v}=\boldsymbol{\theta}} L(\boldsymbol{v}).$

### Theorem

$$I\!\!E\widetilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^*,$$

$$\mathrm{Var}(\widetilde{\boldsymbol{\theta}}) = S\,\mathrm{Var}(\boldsymbol{\varepsilon})S^{\top} = \sigma^2 SS^{\top} = \sigma^2 P(\Upsilon\Upsilon^{\top})^{-1}P^{\top}.$$

Model:

$$\boldsymbol{Y} = \Psi^\top \boldsymbol{\theta}^* + \Phi^\top \boldsymbol{\eta}^* + \boldsymbol{\varepsilon} \qquad I\!\!E \boldsymbol{\varepsilon} = 0, \ \mathrm{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n.$$

### Theorem

*The profile MLE $\widetilde{\boldsymbol{\theta}}$ reads as*

$$\widetilde{\boldsymbol{\theta}} = \big(\breve{\Psi}\breve{\Psi}^\top\big)^{-1} \breve{\Psi} \boldsymbol{Y},$$

$$\breve{\Psi} = \Psi - \Psi \Pi_{\boldsymbol{\eta}} = \Psi - \Psi \Phi^\top \big(\Phi \Phi^\top\big)^{-1} \Phi.$$

Model:

$$\boldsymbol{Y} = \Upsilon^\top \boldsymbol{\upsilon}^* + \boldsymbol{\varepsilon} = \Psi^\top \boldsymbol{\theta}^* + \Phi^\top \boldsymbol{\eta}^* + \boldsymbol{\varepsilon} \qquad I\!\!E\boldsymbol{\varepsilon} = 0, \ \mathrm{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n.$$

**Theorem (Gauss-Markov)**

1. $\widetilde{\boldsymbol{\theta}} = S\boldsymbol{Y}$ with $S = P(\Upsilon\Upsilon^\top)^{-1}\Upsilon$ fulfills

$$I\!\!E\widetilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^* = P\boldsymbol{\upsilon}^*,$$

$$I\!\!E\|\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|^2 = \mathrm{Var}(\widetilde{\boldsymbol{\theta}}) = \sigma^2 P(\Upsilon\Upsilon^\top)^{-1}P^\top = \sigma^2(\breve{\Psi}\breve{\Psi}^\top)^{-1},$$

$$\breve{\Psi} = \Psi - \Psi\Pi_{\boldsymbol{\eta}}$$

$$\Pi_{\boldsymbol{\eta}} = \Phi^\top(\Phi\Phi^\top)^{-1}\Phi.$$

2. This risk is minimal in the class of all unbiased linear estimates of $\boldsymbol{\theta}^*$.

Model:

$$\boldsymbol{Y} = \Psi^\top \boldsymbol{\theta}^* + \Phi^\top \boldsymbol{\eta}^* + \boldsymbol{\varepsilon} \qquad I\!\!E \boldsymbol{\varepsilon} = 0, \ \mathrm{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n.$$

Define

$$\breve{D}_0^2 = \sigma^{-2} \breve{\Psi} \breve{\Psi}^\top, \qquad \breve{\Psi} = \Psi - \Psi \Pi_{\boldsymbol{\eta}}.$$

---

**Theorem**

*Let the matrix $\breve{D}_0^2$ be non-degenerated. It holds*

$$2\{\breve{L}(\widetilde{\boldsymbol{\theta}}) - \breve{L}(\boldsymbol{\theta}^*)\} \ = \ \|\breve{D}_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 = \|\breve{\boldsymbol{\xi}}\|^2,$$

$$\breve{\boldsymbol{\xi}} \ = \ \breve{D}_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*), \qquad I\!\!E \breve{\boldsymbol{\xi}} = 0, \ \mathrm{Var}(\breve{\boldsymbol{\xi}}) = I_p.$$

*If $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$, then $\breve{\boldsymbol{\xi}}$ is standard normal in $I\!\!R^p$ and*

$$2\{\breve{L}(\widetilde{\boldsymbol{\theta}}) - \breve{L}(\boldsymbol{\theta}^*)\} \sim \chi_p^2.$$

SPA

$$\boldsymbol{Y} \sim I\!\!P \in \big(I\!\!P_{\boldsymbol{\theta},\boldsymbol{\eta}},\, \boldsymbol{\theta} \in \Theta, \boldsymbol{\eta} \in H\big)$$

Log-likelihood:

$$L(\boldsymbol{\theta},\boldsymbol{\eta}) \,=\, \frac{dI\!\!P_{\boldsymbol{\theta},\boldsymbol{\eta}}}{d\boldsymbol{\mu}_0}(\boldsymbol{Y})$$

Profile MLE:

$$\widetilde{\boldsymbol{\theta}} \,=\, \operatorname*{argmax}_{\boldsymbol{\theta}} \max_{\boldsymbol{\eta}} L(\boldsymbol{\theta},\boldsymbol{\eta}) = \operatorname*{argmax}_{\boldsymbol{\theta}} \breve{L}(\boldsymbol{\theta}),$$

$$\breve{L}(\boldsymbol{\theta}) \,=\, \max_{\boldsymbol{\eta}} L(\boldsymbol{\theta},\boldsymbol{\eta}).$$

$\boldsymbol{\upsilon}$ -setup: $\boldsymbol{\upsilon} = (\boldsymbol{\theta},\boldsymbol{\eta})\,,\ L(\boldsymbol{\upsilon}) = L(\boldsymbol{\theta},\boldsymbol{\eta})\,,$

$$\widetilde{\boldsymbol{\upsilon}} = \operatorname*{argmax}_{\boldsymbol{\upsilon}} L(\boldsymbol{\upsilon}), \qquad \widetilde{\boldsymbol{\theta}} = P\widetilde{\boldsymbol{\upsilon}}$$

For $L(\boldsymbol{\upsilon}) = L(\boldsymbol{\theta}, \boldsymbol{\eta})$, define

$$\boldsymbol{\upsilon}^* \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\upsilon} \in \Upsilon} I\!\!E L(\boldsymbol{\upsilon}),$$

$$\boldsymbol{\theta}^* \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta}} \max_{\boldsymbol{\eta}} I\!\!E L(\boldsymbol{\theta}, \boldsymbol{\eta}) = P\boldsymbol{\upsilon}^*.$$

Also

$$\mathcal{D}_0^2 \stackrel{\text{def}}{=} -\nabla^2 I\!\!E L(\boldsymbol{\upsilon}^*),$$

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} \mathcal{D}_0^{-1} \nabla L(\boldsymbol{\upsilon}^*),$$

$$\mathcal{V}_0^2 = \operatorname{Var}\{\nabla L(\boldsymbol{\upsilon}^*)\} \quad (= \mathcal{D}_0^2 \ \text{ under PA})$$

and

$$\Upsilon_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{\upsilon} \colon \|\mathcal{D}_0(\boldsymbol{\upsilon} - \boldsymbol{\upsilon}^*)\| \le \mathbf{r}\}.$$

- Concentration and large deviations: fix $\mathbf{r}_0$ ensuring

$$\mathbb{P}\big(\widetilde{\boldsymbol{v}} \notin \Upsilon_\circ(\mathbf{r}_0)\big) \leq \mathrm{e}^{-\mathbf{x}},$$

  where $\quad \Upsilon_\circ(\mathbf{r}) \stackrel{\mathrm{def}}{=} \big\{\boldsymbol{\theta}\colon \|\mathcal{D}_0(\boldsymbol{v} - \boldsymbol{v}^*)\| \leq \mathbf{r}\big\}.$

- Local quadratic approximation of the expected log-likelihood:

$$\sup_{\boldsymbol{v}\in\Upsilon_\circ(\mathbf{r})} \frac{2\mathbb{E}L(\boldsymbol{v}^*) - 2\mathbb{E}L(\boldsymbol{v})}{\|\mathcal{D}_0(\boldsymbol{v} - \boldsymbol{v}^*)\|^2} \leq \delta(\mathbf{r}).$$

- Local linear approximation of the stochastic component: on $\Omega(\mathbf{x})$, for $\zeta(\boldsymbol{v}) \stackrel{\mathrm{def}}{=} L(\boldsymbol{v}) - \mathbb{E}L(\boldsymbol{v})$

$$\sup_{\boldsymbol{v}\in\Upsilon_\circ(\mathbf{r})} \big|\mathcal{D}_0^{-1}\big\{\nabla\zeta(\boldsymbol{v}) - \nabla\zeta(\boldsymbol{v}^*)\big\}\big| \leq \varrho(\mathbf{r},\mathbf{x}).$$

- Overall error of the Fisher expansion $\mathbf{r}_0\big\{\delta(\mathbf{r}_0) + \varrho(\mathbf{r}_0,\mathbf{x})\big\}$, of the Wilks $\mathbf{r}_0^2\big\{\delta(\mathbf{r}_0) + \varrho(\mathbf{r}_0,\mathbf{x})\big\}$.

$$\widetilde{\boldsymbol{v}} \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{v} \in \Upsilon} L(\boldsymbol{v}), \qquad \boldsymbol{v}^* \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{v} \in \Upsilon} I\!\!EL(\boldsymbol{v}),$$

$$\mathcal{D}_0^2 \stackrel{\text{def}}{=} -\nabla^2 I\!\!EL(\boldsymbol{v}^*), \qquad \boldsymbol{\xi} \stackrel{\text{def}}{=} \mathcal{D}_0^{-1} \nabla L(\boldsymbol{v}^*).$$

---

**Theorem**

*On a set* $\Omega(\mathbf{x})$ *with* $I\!\!P\big(\Omega(\mathbf{x})\big) \geq 1 - \mathtt{C}\mathrm{e}^{-\mathbf{x}}$

$$\big\| \mathcal{D}_0(\widetilde{\boldsymbol{v}} - \boldsymbol{v}^*) - \boldsymbol{\xi} \big\| \leq \Diamond(\mathbf{r}_0, \mathbf{x}),$$

$$\Big| L(\widetilde{\boldsymbol{v}}) - L(\boldsymbol{v}^*) - \frac{\|\boldsymbol{\xi}\|^2}{2} \Big| \leq \Delta(\mathbf{r}_0, \mathbf{x}).$$

*Here* $\Diamond(\mathbf{r}_0, \mathbf{x})$ *and* $\Delta(\mathbf{r}_0, \mathbf{x})$ *are explicit error terms.*
*The vector* $\boldsymbol{\xi}$ *fulfills*

$$I\!\!P(\|\boldsymbol{\xi}\| \geq z(B, \mathbf{x})) \leq 2\mathrm{e}^{-\mathbf{x}},$$

*where* $B = \mathrm{Var}(\boldsymbol{\xi}) = \mathcal{D}_0^{-1} \mathcal{V}_0^2 \mathcal{D}_0^{-1}$, *so that* $z^2(B, \mathbf{x}) \asymp p^* + \mathbf{x}$.

Problems: the value of $\|\boldsymbol{\xi}\|^2$ is of order of the full dimension $p^*$.

Corollaries for $\widetilde{\boldsymbol{\theta}} = P\widetilde{\boldsymbol{\upsilon}}$ ?

Consider the block representation:

$$
\mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A \\ A^\top & H_0^2 \end{pmatrix}, \qquad \nabla = \nabla L(\boldsymbol{\upsilon}^*) = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*) \\ \nabla_{\boldsymbol{\eta}} L(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*) \end{pmatrix} = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} \\ \nabla_{\boldsymbol{\eta}} \end{pmatrix},
$$

Define $\breve{D}_0^{-2}$ as the left upper block of $\mathcal{D}_0^{-2}$:

$$
\breve{D}_0^2 = D_0^2 - A H_0^{-2} A^\top
$$

and

$$
\breve{\boldsymbol{\xi}} \stackrel{\text{def}}{=} \breve{D}_0^{-1} \big( \nabla_{\boldsymbol{\theta}} - A H_0^{-2} \nabla_{\boldsymbol{\eta}} \big)
$$

$$\mathcal{D}_0^2 = \begin{pmatrix} D_0^2 & A \\ A^\top & H_0^2 \end{pmatrix}, \qquad \nabla = \nabla L(\boldsymbol{v}^*) = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} \\ \nabla_{\boldsymbol{\eta}} \end{pmatrix},$$

$$\breve{D}_0^2 = D_0^2 - A H_0^{-2} A^\top \qquad \breve{\boldsymbol{\xi}} \stackrel{\text{def}}{=} \breve{D}_0^{-1} \breve{\nabla}_{\boldsymbol{\theta}} = \breve{D}_0^{-1}\big(\nabla_{\boldsymbol{\theta}} - A H_0^{-2} \nabla_{\boldsymbol{\eta}}\big)$$

**Theorem**

*On a set $\Omega(\mathbf{x})$ with $I\!P\big(\Omega(\mathbf{x})\big) \geq 1 - \mathtt{C}\,e^{-\mathbf{x}}$*

$$\big\| \breve{D}_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \breve{\boldsymbol{\xi}} \big\| \leq \Diamond(\mathbf{r}_0, \mathbf{x}),$$

$$\big| \breve{L}(\widetilde{\boldsymbol{\theta}}) - \breve{L}(\boldsymbol{\theta}^*) - \frac{\|\breve{\boldsymbol{\xi}}\|^2}{2} \big| \leq \Delta(\mathbf{r}_0, \mathbf{x}) \leq \mathtt{C}\, p \,\Diamond(\mathbf{r}_0, \mathbf{x}).$$

*Here $\Diamond(\mathbf{x})$ and $\Delta(\mathbf{x})$ are explicit error terms. The vector $\breve{\boldsymbol{\xi}}$ fulfills*

$$I\!P(\|\breve{\boldsymbol{\xi}}\| \geq z(\breve{B}, \mathbf{x})) \leq 2e^{-\mathbf{x}},$$

*where $\breve{B} = \operatorname{Var}(\breve{\boldsymbol{\xi}}) = \breve{D}_0^{-1} \operatorname{Var}(\breve{\nabla}) \breve{D}_0^{-1}$, so that $z^2(\breve{B}, \mathbf{x}) \asymp p + \mathbf{x}$.*

**Concentration and large deviation for the PMLE $\widetilde{\theta}$**

Steps:

- Concentration of $\widetilde{\boldsymbol{v}}$ on $\Upsilon_\circ(\mathbf{r}_0)$ for $\mathbf{r}_0^2 \asymp p^* + \mathbf{x}$;
- Full dimensional Fisher expansion: on $\Omega(\mathbf{x})$

$$\big\| \mathcal{D}_0(\widetilde{\boldsymbol{v}} - \boldsymbol{v}^*) - \boldsymbol{\xi} \big\| \leq \diamondsuit(\mathbf{r}_0, \mathbf{x});$$

- Fisher expansion for $\widetilde{\boldsymbol{\theta}}$: on $\Omega(\mathbf{x})$

$$\big\| \breve{D}_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \breve{\boldsymbol{\xi}} \big\| \leq \diamondsuit(\mathbf{r}_0, \mathbf{x});$$

- A deviation bound

$$I\!P(\|\breve{\boldsymbol{\xi}}\| \geq z(\breve{B}, \mathbf{x})) \leq 2\mathrm{e}^{-\mathbf{x}}$$

Imply concentration of $\widetilde{\boldsymbol{\theta}}$ on $\Theta_0(\breve{\mathbf{r}}_0)$ for $\breve{\mathbf{r}}_0 = z(\breve{B}, \mathbf{x}) + \diamondsuit(\mathbf{r}_0, \mathbf{x})$:

$$I\!P\Big\{ \big\| \breve{D}_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \big\| \geq z(\breve{B}, \mathbf{x}) + \diamondsuit(\mathbf{r}_0, \mathbf{x}) \Big\} \leq 3\mathrm{e}^{-\mathbf{x}}.$$

# Outline

Let $p = p_n \to \infty$. We know

$$\diamondsuit_n(\mathbf{x}) \leq \mathtt{C} \sqrt{\frac{(p_n + \mathbf{x})^2}{n}}, \qquad \Delta_n(\mathbf{x}) \leq \mathtt{C} \sqrt{\frac{(p_n + \mathbf{x})^3}{n}}, \qquad \|\boldsymbol{\xi}_n\|^2 \leq p_n + \mathtt{C}\mathbf{x}.$$

■ $p_n/n \to 0$ : Consistency:

$$\|\sqrt{\mathbb{F}_{\boldsymbol{\theta}^*}}(\widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*)\| = n^{-1/2}\{\|\boldsymbol{\xi}_n\| \pm \diamondsuit_n(\mathbf{x})\} \leq \sqrt{\frac{p_n + \mathtt{C}\mathbf{x}}{n}} \pm \mathtt{C}\frac{p_n + \mathbf{x}}{n}$$

■ $p_n^2/n \to 0$ – Fisher expansion, root-$n$ normality;

$$\sqrt{n\mathbb{F}_{\boldsymbol{\theta}^*}}(\widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) = \boldsymbol{\xi}_n \pm \diamondsuit_n(\mathbf{x}), \qquad \text{expansion of the MLE}$$

$$\sqrt{2L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)} = \|\boldsymbol{\xi}_n\| \pm 3\diamondsuit_n(\mathbf{x}), \qquad \text{square-root excess}$$

$$p_n^{-1/2}L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = p_n^{-1/2}\|\boldsymbol{\xi}_n\|^2/2 \pm \mathtt{C}\diamondsuit_n(\mathbf{x}), \qquad \text{likelihood ratio tests, model selection}$$

■ $p_n^3/n \to 0$ – Wilks approximation, BvM Theorem.

## Penalization

Let $\mathrm{pen}(\boldsymbol{\theta})$ be a penalty function on $\Theta$.

Large $\mathrm{pen}(\boldsymbol{\theta}) \iff$ rough $\boldsymbol{\theta}$.

Small $\mathrm{pen}(\boldsymbol{\theta}) \iff$ smooth $\boldsymbol{\theta}$.

Structural assumption – the true value $\boldsymbol{\theta}^*$ is smooth – $\mathrm{pen}(\boldsymbol{\theta}_0)$ is (relatively) small.

A penalized (quasi) MLE approach leads to maximizing the penalized log-likelihood:

$$\widetilde{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \big\{ L(\boldsymbol{\theta}) - \mathrm{pen}(\boldsymbol{\theta}) \big\}.$$

New target:

$$\boldsymbol{\theta}^*_{\mathrm{pen}} = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \big\{ I\!\!E L(\boldsymbol{\theta}) - \mathrm{pen}(\boldsymbol{\theta}) \big\}.$$

In general, $\boldsymbol{\theta}^* \neq \boldsymbol{\theta}^*_{\mathrm{pen}}$: "modeling bias" issue.

**Quadratic penalization**

Important special case – a quadratic penalty $\mathrm{pen}(\boldsymbol{\theta}) = \|G\boldsymbol{\theta}\|^2/2$ for a given symmetric matrix $G^2$. Denote

$$L_G(\boldsymbol{\theta}) \stackrel{\mathrm{def}}{=} L(\boldsymbol{\theta}) - \|G\boldsymbol{\theta}\|^2/2,$$

$$\widetilde{\boldsymbol{\theta}}_G \stackrel{\mathrm{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L_G(\boldsymbol{\theta}).$$

The use of a penalty changes the target of estimation which is now defined as

$$\boldsymbol{\theta}_G^* \stackrel{\mathrm{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} I\!E L_G(\boldsymbol{\theta}).$$

In general $\boldsymbol{\theta}^* \neq \boldsymbol{\theta}_G^*$.

The modeling bias can be measured by $\|G\boldsymbol{\theta}^*\|^2$, yielding the "bias-variance" trade-off:

$$I\!E\|\boldsymbol{\xi}_G\|^2 \asymp \|G\boldsymbol{\theta}^*\|^2$$

Let $V_0^2 = \mathrm{Var}\{\nabla L(\boldsymbol{\theta}_G^*)\}$.

Typically $V_0^2$ measures the variability of the process $L(\cdot)$ and $L_G(\cdot)$.

Let also $D_G^2$ be a penalized information matrix

$$D_G^2 = -\nabla^2 I\!\!E L_G(\boldsymbol{\theta}_G^*) = D_0^2 + G^2$$

with $D_0^2 = -\nabla^2 I\!\!E L(\boldsymbol{\theta}_G^*)$.

The effective dimension $\mathrm{p}_G$ is defined as the trace of the matrix $B_G \stackrel{\text{def}}{=} D_G^{-1} V_0^2 D_G^{-1}$:

$$\mathrm{p}_G \stackrel{\text{def}}{=} \mathrm{tr}\big(B_G\big) = I\!\!E \|\boldsymbol{\xi}_G\|^2$$

for $\boldsymbol{\xi}_G = D_G^{-1} \nabla L(\boldsymbol{\theta}_G^*)$.

Let

$$V_0^2 = D_0^2 = \sigma^2 I_p,$$
$$G^2 = \text{diag}\{g_1^2 \geq g_2^2 \geq \ldots g_p^2\}$$

Then

$$D_G^2 = D_0^2 + G^2 = \text{diag}\{\sigma^2 + g_1^2, \ldots, \sigma^2 + g_p^2\},$$
$$B_G = \text{diag}\{(1 + \sigma^{-2} g_1^2)^{-1}, \ldots, (1 + \sigma^{-2} g_p^2)^{-1}\}.$$

$G$ is of a block structure: $G = \operatorname{diag}\{0, G_1\}$.

The first block of dimension $p_0$ corresponds to the unconstrained part of the parameter vector

the second block of dimension $p_1$ corresponds to the low energy component.

Assume for simplicity that $G_1 = g I_{p_1}$. Then

$$\mathfrak{p}_G = \operatorname{tr} B_G = p_0 + p_1 / \left( 1 + \sigma^{-2} g^2 \right).$$

The impact of $G_1$ in the effective dimension is inessential if $g^2 / \sigma^2 \gg p_1 / p_0$.

For $\beta > 1/2$,

$$G^2 = \mathrm{diag}\{g_1^2, \ldots, g_p^2\}$$

$$g_j = Lj^\beta$$

The value $\beta$ is usually considered as the Sobolev smoothness parameter.

It holds

$$\mathtt{p}_G = \sum_{j=1}^{p} \frac{1}{1 + L^2 j^{2\beta}/\sigma^2} \, .$$

Define also the index $\mathtt{p}_e$ as the largest $j$ satisfying $Lj^\beta \leq \sigma$.

$\beta > 1/2$ yields $\mathtt{p}_G \leq \mathtt{C}(\beta)\mathtt{p}_e$ for some constant $\mathtt{C}(\beta)$ depending on $\beta$ only.

$$\widetilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}), \qquad \boldsymbol{\theta}^* \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} I\!\!E L(\boldsymbol{\theta})$$

**Theorem**

*On a set* $\Omega(\mathbf{x})$ *with* $I\!\!P\big(\Omega(\mathbf{x})\big) \geq 1 - \mathsf{C}e^{-\mathbf{x}}$

$$\big\| D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi} \big\| \ \leq \ \Diamond(\mathbf{x}),$$

$$\big| L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) - \frac{\|\boldsymbol{\xi}\|^2}{2} \big| \ \leq \ \Delta(\mathbf{x})$$

*with*

$$D_0^2 \stackrel{\text{def}}{=} -\nabla^2 I\!\!E L(\boldsymbol{\theta}^*), \qquad \boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla L(\boldsymbol{\theta}^*).$$

$$\widetilde{\boldsymbol{\theta}}_G \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L_G(\boldsymbol{\theta}), \qquad \boldsymbol{\theta}_G^* \stackrel{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E} L_G(\boldsymbol{\theta})$$

**Theorem**

On a set $\Omega(\mathtt{x})$ with $\mathbb{P}\big(\Omega(\mathtt{x})\big) \geq 1 - \mathtt{C}\mathrm{e}^{-\mathtt{x}}$

$$\big\| D_G(\widetilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) - \boldsymbol{\xi}_G \big\| \leq \Diamond_G(\mathtt{x}),$$

$$\big| L_G(\widetilde{\boldsymbol{\theta}}_G) - L_G(\boldsymbol{\theta}_G^*) - \frac{\|\boldsymbol{\xi}_G\|^2}{2} \big| \leq \Delta_G(\mathtt{x})$$

with

$$D_G^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E} L_G(\boldsymbol{\theta}_G^*) = -\nabla^2 \mathbb{E} L(\boldsymbol{\theta}_G^*) + G^2,$$

$$\boldsymbol{\xi}_G \stackrel{\text{def}}{=} D_G^{-1} \nabla L_G(\boldsymbol{\theta}^*).$$

$(\mathcal{L})$  *For each* $\mathbf{r}$, *there exists* $\mathbf{b}(\mathbf{r}) > 0$ *such that* $\mathbf{rb}(\mathbf{r}) \to \infty$ *as* $\mathbf{r} \to \infty$ *and*

$$\frac{-2\mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2} \geq \mathbf{b}(\mathbf{r}), \quad \forall \boldsymbol{\theta} \in \Theta_0(\mathbf{r}) = \{\boldsymbol{\theta} \colon \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}.$$

---

**Theorem**

*Suppose* $(ED_0)$ *and* $(ED_2)$, $(\mathcal{L}_0)$, $(\mathcal{L})$, *and* $(\mathcal{I})$. *Let* $\mathbf{b}(\mathbf{r})$ *in* $(\mathcal{L})$ *satisfy*

$$\mathbf{b}(\mathbf{r})\,\mathbf{r} \geq 2z(B, \mathbf{x}) + 2\varrho(\mathbf{r}, \mathbf{x}), \quad \mathbf{r} > \mathbf{r}_0,$$

*where* $B = D_0^{-1} V_0^2 D_0$ *and*

$$\varrho(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 6\nu_0\, z_{\mathbb{H}}\big(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)\big)\, \omega. \tag{13}$$

*Then*

$$\mathbb{P}\big(\widetilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}_0)\big) \leq 3\mathrm{e}^{-\mathbf{x}}.$$

$(\mathcal{L}G)$   *For each* $\mathbf{r}$ *, there exists* $\mathbf{b}_G(\mathbf{r}) > 0$ *such that* $\mathbf{r}\mathbf{b}_G(\mathbf{r}) \to \infty$ *as* $\mathbf{r} \to \infty$ *and*

$$\frac{-2I\!\!E L_G(\boldsymbol{\theta}, \boldsymbol{\theta}_G^*)}{\|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\|^2} \geq \mathbf{b}_G(\mathbf{r}), \quad \forall \boldsymbol{\theta} \in \Theta_{0,G}(\mathbf{r}) = \big\{ \boldsymbol{\theta} \colon \|D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*)\| \leq \mathbf{r} \big\}.$$

**Theorem**

*Let* $\mathbf{b}_G(\mathbf{r})$ *in* $(\mathcal{L}G)$ *satisfy*

$$\mathbf{b}_G(\mathbf{r})\,\mathbf{r} \geq 2z(B_G, \mathbf{x}) + 2\varrho(\mathbf{r}, \mathbf{x}), \quad \mathbf{r} > \mathbf{r}_0,$$

*where* $B_G = D_G^{-1} V_0^2 D_G$

$$\varrho(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 6\nu_0\, z_{\mathbb{H}}\big(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)\big)\, \omega. \tag{14}$$

*Then*

$$I\!\!P\big(\widetilde{\boldsymbol{\theta}}_G \notin \Theta_{0,G}(\mathbf{r}_0)\big) \leq 3\mathrm{e}^{-\mathbf{x}}.$$

Let a vector process $\mathcal{Y}(\boldsymbol{v})$ fulfill on $\Upsilon_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{v} \colon \|\boldsymbol{v}\| \le \mathbf{r}\}$

$$\sup_{\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in \mathbb{R}^p \colon \|\boldsymbol{\gamma}_1\|=\|\boldsymbol{\gamma}_2\|=1} \log I\!\!E \exp\left\{\lambda \boldsymbol{\gamma}_1^\top \nabla \mathcal{Y}(\boldsymbol{v}) \boldsymbol{\gamma}_2\right\} \le \frac{\nu_0^2 \lambda^2}{2}.$$

---

**Theorem**

*Suppose* $(ED_2)$ *. It holds on a random set* $\Omega(\mathbf{r}, \mathbf{x})$

$$\sup_{\boldsymbol{v} \in \Upsilon_\circ(\mathbf{r})} \|\mathcal{Y}(\boldsymbol{v})\| \le 6\nu_0 \, z_{\mathbb{H}}(\mathbf{x}) \, \mathbf{r},$$

*where the function* $z_{\mathbb{H}}(\mathbf{x})$ *is given by*

$$z_{\mathbb{H}}(\mathbf{x}) = \mathbb{H}_1 + \sqrt{2\mathbf{x}} + \mathbf{g}^{-1}(\mathbf{g}^{-2}\mathbf{x} + 1)\mathbb{H}_2,$$

*with* $\mathbb{H}_2 = 4p$ *and* $\mathbb{H}_1 = 2p^{1/2}$ *.*

---

**A bound for the norm of a vector stochastic process "penalized"**

Let a vector process $\mathcal{Y}(\boldsymbol{v})$ fulfill on $\Upsilon_{\circ}(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{v} \colon \|B_G^{-1/2}\boldsymbol{v}\| \leq \mathbf{r}\}$

$$\sup_{\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in \mathbb{R}^p \colon \|\boldsymbol{\gamma}_1\|=\|\boldsymbol{\gamma}_2\|=1} \log \mathbb{E} \exp\left\{\lambda \boldsymbol{\gamma}_1^\top \nabla \mathcal{Y}(\boldsymbol{v}) \boldsymbol{\gamma}_2\right\} \leq \frac{\nu_0^2 \lambda^2}{2}.$$

**Theorem**

*Suppose* $(ED_2)$ *. It holds on a random set* $\Omega(\mathbf{r}, \mathbf{x})$

$$\sup_{\boldsymbol{v} \in \Upsilon_{\circ}(\mathbf{r})} \|B_G^{1/2}\mathcal{Y}(\boldsymbol{v})\| \leq 6\nu_0 \, z_{\mathbb{H}}(\mathbf{x}) \, \mathbf{r},$$

*where the function* $z_{\mathbb{H}}(\mathbf{x})$ *is given by*

$$z_{\mathbb{H}}(\mathbf{x}) = \mathbb{H}_1 + \sqrt{2\mathbf{x}} + \mathbf{g}^{-1}(\mathbf{g}^{-2}\mathbf{x} + 1)\mathbb{H}_2,$$

*with*

$$\mathbb{H}_1 = \mathbb{H}_1(B_G) = 1 + 2\sqrt{\text{tr}(B_G \log(B_G))}, \quad \mathbb{H}_2 = \mathbb{H}_2(B) = 1 + \frac{8}{3}\text{tr}(B_G^{1/2}).$$

On $\Omega(\mathbf{r}, \mathbf{x})$, for each $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$

$$\left\| D_0^{-1} \left\{ \nabla I\!E L(\boldsymbol{\theta}) - \nabla I\!E L(\boldsymbol{\theta}^*) \right\} + D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \leq \delta(\mathbf{r})\mathbf{r},$$

$$\left\| D_0^{-1} \left\{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*) \right\} \right\| \leq 6\nu_0 \, z_{\mathbb{H}}(\mathbf{x}) \, \omega \, \mathbf{r}$$

---

**Theorem**

*Suppose $(\mathcal{L}_0)$ and $(ED_2)$ on $\Theta_0(\mathbf{r})$ for a fixed $\mathbf{r}$. Then on $\Omega(\mathbf{r}, \mathbf{x})$*

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\| D_0^{-1} \left\{ \nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^*) \right\} + D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \leq \Diamond(\mathbf{r}, \mathbf{x}),$$

*where*

$$\Diamond(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \left\{ \delta(\mathbf{r}) + 6\nu_0 \, z_{\mathbb{H}}(\mathbf{x}) \, \omega \right\} \mathbf{r}.$$

The dimension $p$ enters only via the entropy $\mathbb{H}$ in $z_{\mathbb{H}}(\mathbf{x})$.

**Local linear approximation of the gradient "penalized"**

On $\Omega(\mathbf{r}, \mathbf{x})$, for each $\boldsymbol{\theta} \in \Theta_{0,G}(\mathbf{r})$

$$\left\| D_G^{-1}\left\{ \nabla I\!\!E L_G(\boldsymbol{\theta}) - \nabla I\!\!E L_G(\boldsymbol{\theta}_G^*) \right\} + D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*) \right\| \leq \delta_G(\mathbf{r})\mathbf{r},$$

$$\left\| D_G^{-1}\left\{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}_G^*) \right\} \right\| \leq 6\nu_0\, z_{\mathbb{H}}(\mathbf{x})\, \omega\, \mathbf{r}$$

---

**Theorem**

*Suppose* $(\mathcal{L}_0 G)$ *and* $(ED_2 G)$ *on* $\Theta_{0,G}(\mathbf{r})$ *for a fixed* $\mathbf{r}$. *Then on* $\Omega(\mathbf{r}, \mathbf{x})$

$$\sup_{\boldsymbol{\theta} \in \Theta_{0,G}(\mathbf{r})} \left\| D_G^{-1}\left\{ \nabla L_G(\boldsymbol{\theta}) - \nabla L_G(\boldsymbol{\theta}_G^*) \right\} + D_G(\boldsymbol{\theta} - \boldsymbol{\theta}_G^*) \right\| \leq \Diamond_G(\mathbf{r}, \mathbf{x}),$$

*where*

$$\Diamond_G(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \left\{ \delta_G(\mathbf{r}) + 6\nu_0\, z_{\mathbb{H}}(\mathbf{x})\, \omega \right\}\mathbf{r}.$$

The effective dimension $\mathrm{p}_G$ enters only via the entropy $\mathbb{H}$ in $z_{\mathbb{H}}(\mathbf{x})$.

Let $p = p_n \to \infty$. We know

$$\Diamond_n(\mathtt{x}) \leq \mathtt{C}\sqrt{\frac{(p_n + \mathtt{x})^2}{n}}, \qquad \Delta_n(\mathtt{x}) \leq \mathtt{C}\sqrt{\frac{(p_n + \mathtt{x})^3}{n}}, \qquad \|\boldsymbol{\xi}_n\|^2 \leq p_n + \mathtt{C}\mathtt{x}.$$

■ $p_n/n \to 0$ : Consistency:

$$\|\sqrt{\mathbb{F}_{\boldsymbol{\theta}^*}}(\widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*)\| = n^{-1/2}\{\|\boldsymbol{\xi}_n\| \pm \Diamond_n(\mathtt{x})\} \leq \mathtt{C}\sqrt{\frac{p_n + \mathtt{x}}{n}} \pm \mathtt{C}\frac{p_n + \mathtt{x}}{n}$$

■ $p_n^2/n \to 0$ – Fisher expansion, root-$n$ normality;

$$\sqrt{n\mathbb{F}_{\boldsymbol{\theta}^*}}(\widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) = \boldsymbol{\xi}_n \pm \Diamond_n(\mathtt{x}), \qquad \text{Expansion of the MLE}$$

$$\sqrt{2L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)} = \|\boldsymbol{\xi}_n\| \pm 3\Diamond_n(\mathtt{x}), \qquad \text{square-root maximum likelihood}$$

$$p_n^{-1/2}L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = p_n^{-1/2}\|\boldsymbol{\xi}_n\|^2/2 \pm \mathtt{C}\Diamond_n(\mathtt{x}), \qquad \text{likelihood ratio tests, model selection}$$

■ $p_n^3/n \to 0$ – Wilks approximation, BvM Theorem.

Let $p = p_n \to \infty$. We know

$$\diamondsuit_G(\mathbf{x}) \leq \mathtt{C} \sqrt{\frac{(\mathtt{p}_G + \mathbf{x})^2}{n}}, \qquad \Delta_G(\mathbf{x}) \leq \mathtt{C} \sqrt{\frac{(\mathtt{p}_G + \mathbf{x})^3}{n}}, \qquad \|\boldsymbol{\xi}_G\|^2 \leq \mathtt{p}_G + \mathtt{C}\,\mathbf{x}.$$

- $\mathtt{p}_G/n \to 0$: Consistency: with $\mathbb{F}_G = \mathbb{F}_{\boldsymbol{\theta}_G^*} + n^{-1}G^2$

$$\|\sqrt{\mathbb{F}_G}(\widetilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*)\| = n^{-1/2}\{\|\boldsymbol{\xi}_G\| \pm \diamondsuit_G(\mathbf{x})\} \leq \mathtt{C}\,\sqrt{\frac{\mathtt{p}_G + \mathbf{x}}{n}} \pm \mathtt{C}\,\frac{\mathtt{p}_G + \mathbf{x}}{n}$$

- $\mathtt{p}_G^2/n \to 0$ – Fisher expansion, root-$n$ normality;

$$\sqrt{n\mathbb{F}_G}(\widetilde{\boldsymbol{\theta}}_G - \boldsymbol{\theta}_G^*) = \boldsymbol{\xi}_G \pm \diamondsuit_G(\mathbf{x}), \qquad \text{Expansion of the MLE}$$

$$\sqrt{2L_G(\widetilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*)} = \|\boldsymbol{\xi}_G\| \pm 3\diamondsuit_G(\mathbf{x}), \qquad \text{square-root maximum likelihood}$$

$$\mathtt{p}_G^{-1/2}L_G(\widetilde{\boldsymbol{\theta}}_G, \boldsymbol{\theta}_G^*) = \mathtt{p}_G^{-1/2}\|\boldsymbol{\xi}_G\|^2/2 \pm \mathtt{C}\,\diamondsuit_G(\mathbf{x}), \qquad \text{likelihood ratio tests, model selection}$$

- $\mathtt{p}_G^3/n \to 0$ – Wilks approximation, BvM Theorem.

# Outline

The $1 - \alpha$ confidence set for $\boldsymbol{\theta}^*$ :

$$\mathcal{E}(\mathfrak{z}_\alpha) \stackrel{\text{def}}{=} \{\boldsymbol{\theta} \colon L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}) \leq \mathfrak{z}_\alpha\},$$

$$I\!\!P\left(\boldsymbol{\theta}^* \notin \mathcal{E}(\mathfrak{z}_\alpha)\right) \leq \alpha.$$

For the known $L(\boldsymbol{\theta})$ and $\alpha$ the set is determined by the critical value $\mathfrak{z}_\alpha$, the $1 - \alpha$ quantile of the excess $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)$.

For $L(\boldsymbol{\theta}) = -\|\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta}\|^2/2$, $\mathcal{E}(\mathfrak{z})$ is an ellipsoid:

$$\mathcal{E}(\mathfrak{z}) = \{\boldsymbol{\theta} \colon \|\Psi^\top (\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}})\|^2 \leq 2\mathfrak{z}\}.$$

▶ Under PA, in the asymptotic setup, $\mathfrak{z}_\alpha$ is close to $1 - \alpha$ quantiles of $\chi_p^2$ due to the Wilks phenomenon:

$$L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \approx \|\boldsymbol{\xi}_n\|^2/2, \qquad \boldsymbol{\xi}_n \xrightarrow{w} \mathcal{N}(0, \boldsymbol{I}_p), \quad n \to \infty.$$

▶ But the speed of convergence is slow and under PA-PW the limit distribution is non-pivotal, i.e. depends on $I\!P$.

▶ The non-asymptotic Wilks result cannot help directly, since the deviation bound for $\|\boldsymbol{\xi}\|^2$ is also non-pivotal and is too rough for a sharp confidence set

$$\left| L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2/2 \right| \leq \Delta(\mathbf{x}),$$

$$I\!P\left( \|\boldsymbol{\xi}\|^2 \geq \mathtt{C}(p + 6\mathbf{x}) \right) \leq 2\mathrm{e}^{-\mathbf{x}}.$$

The idea is to mimic the distribution of $L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$ using multiplier bootstrap.

Below $\ell_i(\boldsymbol{\theta})$ is the log-density of $Y_i$ : $\ell_i(\boldsymbol{\theta}) = \log \frac{dP_i(\boldsymbol{\theta})}{d\mu_0}(Y_i)$ and

$$L(\boldsymbol{\theta}) = \sum_{i=1}^{n} \ell_i(\boldsymbol{\theta}).$$

- Take an i.i.d. sample $u_1, \ldots, u_n$ independent of the data $\boldsymbol{Y}$, $I\!\!E(u_i) = \text{Var}(u_i) = 1$ (e.g. $u_i \sim \exp(1)$ or $\mathcal{N}(1,1)$ ).

- Bootstrap the likelihood function:

$$L^{\circ}(\boldsymbol{\theta}) = L^{\circ}(\boldsymbol{\theta}, \boldsymbol{u}) \stackrel{\text{def}}{=} \sum_{i=1}^{n} \ell_i(\boldsymbol{\theta}) \, u_i$$

$^{\circ}$ denotes the conditional probability with the fixed sample $\boldsymbol{Y}$ .

| "$\boldsymbol{Y}$ world" | "bootstrap world" |
|---|---|
| MLE | |
| $\widetilde{\boldsymbol{\theta}} \overset{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$ | $\widetilde{\boldsymbol{\theta}}^{\boldsymbol{\circ}} \overset{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} L^{\boldsymbol{\circ}}(\boldsymbol{\theta})$ |
| target | |
| $\boldsymbol{\theta}^* \overset{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} I\!\!E L(\boldsymbol{\theta})$ | $\widetilde{\boldsymbol{\theta}} \overset{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta}} I\!\!E^{\boldsymbol{\circ}} L^{\boldsymbol{\circ}}(\boldsymbol{\theta})$ |
| likelihood ratio | |
| $L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*)$ | $L^{\boldsymbol{\circ}}(\widetilde{\boldsymbol{\theta}}^{\boldsymbol{\circ}}) - L^{\boldsymbol{\circ}}(\widetilde{\boldsymbol{\theta}})$ |

■ The bootstrap side is fully computable!

■ The true point in bootstrap world is exactly qMLE $\widetilde{\boldsymbol{\theta}}$ .

■ The "bootstrap world" is built inside of the parametric model, which may be wrong.

**Questions to be addressed:**

- Bootstrap consistency in non-asymptotic form

- Error of coverage probability

- Size of the bootstrap-based confidence set

**Key ingredients:**

- Fisher and Wilks expansions in real and bootstrap worlds;

- Closeness of distributions of the of approximating terms $\|\boldsymbol{\xi}\|^2$ and $\|\boldsymbol{\xi}^{\circledast}\|^2$ ;

- Closeness of the local metrics on the parameter space:

$$D_0^2 \approx \mathfrak{D}_0^2 \quad \Leftrightarrow \quad \nabla_{\boldsymbol{\theta}}^2 I\!\!E L(\boldsymbol{\theta}^*) \approx \nabla_{\boldsymbol{\theta}}^2 I\!\!E^{\circledast} L^{\circledast}(\widetilde{\boldsymbol{\theta}});$$

- Use of the truncated moment-generating function to get a sharp bound for

$$\mathcal{L}(\|\boldsymbol{\xi}\|^2) \approx \mathcal{L}(\|\boldsymbol{\xi}^{\circledast}\|^2 \mid \boldsymbol{Y}).$$

---

**Theorem**

*It holds with* $\mathbb{P}^{\bullet}-$ *probability* $\geq 1 - C\,\mathrm{e}^{-\mathbf{x}}$ .

$$\left| L^{\bullet}(\widetilde{\boldsymbol{\theta}}^{\bullet}, \widetilde{\boldsymbol{\theta}}) - \|\boldsymbol{\xi}^{\bullet}\|^2/2 \right| \leq \Delta^{\bullet}(\mathbf{x}),$$

*where the following terms are* $\mathbb{P}-$ *random*

$$\boldsymbol{\xi}^{\bullet} \stackrel{\mathrm{def}}{=} \mathfrak{D}_0^{-1} \nabla_{\boldsymbol{\theta}} L^{\bullet}(\widetilde{\boldsymbol{\theta}}), \qquad \mathfrak{D}_0^2 \stackrel{\mathrm{def}}{=} -\nabla_{\boldsymbol{\theta}}^2 \mathbb{E}^{\bullet} L^{\bullet}(\widetilde{\boldsymbol{\theta}}).$$

Two Wilks results lead to the following scheme:

$$L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \approx \|\boldsymbol{\xi}\|^2/2$$

$$\wr\wr$$

$$L^{\circledast}(\widetilde{\boldsymbol{\theta}}^{\circledast}, \widetilde{\boldsymbol{\theta}}) \approx \|\boldsymbol{\xi}^{\circledast}\|^2/2.$$

The Wilks theorems results are valid on two different probability spaces. The approximation $\approx$ connects two "worlds" in distribution:

$$\mathcal{L}(\|\boldsymbol{\xi}\|^2) \approx \mathcal{L}(\|\boldsymbol{\xi}^{\circledast}\|^2 \,|\, \boldsymbol{Y}).$$

Leading to

$$\mathcal{L}\big\{L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)\big\} \approx \mathcal{L}\big\{L^{\circledast}(\widetilde{\boldsymbol{\theta}}^{\circledast}, \widetilde{\boldsymbol{\theta}}) \,|\, \boldsymbol{Y}\big\}.$$

**Theorem**

*Let the conditions $(ED_2)$, $(ED_3)$ and $(\mathcal{L}_0)$ be fulfilled, then it holds with probability $\geq 1 - 2e^{-\mathtt{x}}$*

$$\sup_{\substack{\boldsymbol{\gamma}_{1,2} \in \mathbb{R}^p, \\ \|\boldsymbol{\gamma}_{1,2}\|=1}} \sup_{\boldsymbol{\theta} \in \Theta_0(r_0)} \left| \boldsymbol{\gamma}_1^\top D_0^{-1} \mathfrak{D}^2(\boldsymbol{\theta}) D_0^{-1} \boldsymbol{\gamma}_2 - 1 \right| \leq C\sqrt{(p+\mathtt{x})^3/n},$$

*where*

$$\mathfrak{D}^2(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\sum_{i=1}^n \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}), \qquad D_0^2 \stackrel{\text{def}}{=} -\sum_{i=1}^n \mathbb{E}\nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}^*).$$

This result implies that on the set $\Omega(\mathtt{x})$ of a dominating probability $1 - C\,e^{\mathtt{x}}$

$$\|D_0^{-1}\mathfrak{D}_0^2 D_0^{-1} - \boldsymbol{I}_p\|_\infty \leq C\sqrt{(p+\mathtt{x})^3/n}.$$

**Lemma**

It holds with $I\!\!P^{\bullet}-$ probability $\geq 1 - 2e^{-x}$

$$\sup_{\boldsymbol{\theta}_1,\boldsymbol{\theta}_2 \in \Theta_0^{\bullet}(\mathbf{r}_0)} \|\boldsymbol{\xi}^{\bullet}(\boldsymbol{\theta}_1) - \boldsymbol{\xi}^{\bullet}(\boldsymbol{\theta}_2)\| \leq C(p+\mathbf{x})/\sqrt{n}.$$

Moreover

$$\left| \|\boldsymbol{\xi}^{\bullet}(\widetilde{\boldsymbol{\theta}})\|^2 - \|\boldsymbol{\xi}^{\bullet}(\boldsymbol{\theta}^*)\|^2 \right| \leq C\sqrt{(p+\mathbf{x})^3/n},$$

where

$$\boldsymbol{\xi}^{\bullet}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathfrak{D}_0^{-1} \left\{ \nabla_{\boldsymbol{\theta}} L^{\bullet}(\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} I\!\!E^{\bullet} L^{\bullet}(\boldsymbol{\theta}) \right\},$$

$$\Theta_0^{\bullet}(\mathbf{r}_0) \stackrel{\text{def}}{=} \left\{ \boldsymbol{\theta} \colon \|\mathfrak{D}_0(\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}})\| \leq \mathbf{r}_0 \right\}.$$

Remind the definition:

Normalized score functions:

$$\boldsymbol{\xi} = D_0^{-1} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*) = D_0^{-1} \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*),$$

$$\boldsymbol{\xi}^{\circ} = \mathfrak{D}_0^{-1} \nabla_{\boldsymbol{\theta}} L^{\circ}(\widetilde{\boldsymbol{\theta}}) = \mathfrak{D}_0^{-1} \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}} \ell_i(\widetilde{\boldsymbol{\theta}})(\boldsymbol{u_i - 1}).$$

Fisher Information matrices

$$D_0^2 = -\sum_{i=1}^{n} I\!\!E \nabla_{\boldsymbol{\theta}}^2 \ell_i(\boldsymbol{\theta}^*) \quad \text{deterministic,}$$

$$\mathfrak{D}_0^2 = -\sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}}^2 \ell_i(\widetilde{\boldsymbol{\theta}}) \quad I\!\!P - \text{random.}$$

Due to the previous results one can make the following substitution: on a set of probability $\geq 1 - \mathtt{C}e^{-\mathtt{x}}$:

$$\mathfrak{D}_0^2 \approx D_0^2, \qquad \|\boldsymbol{\xi}^{\bullet}(\widetilde{\theta})\|^2 \approx \|\boldsymbol{\xi}^{\bullet}(\theta^*)\|^2, \qquad \boldsymbol{\xi}^{\bullet}(\widetilde{\theta}) \approx \boldsymbol{\xi}^{\bullet}(\theta^*).$$

$$\boldsymbol{\xi}^{\bullet}(\widetilde{\theta}) \approx \boldsymbol{\xi}^{\bullet}(\theta^*) \;\approx\; D_0^{-1} \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)(\boldsymbol{u_i - 1}),$$

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*).$$

**Multiplier CLT**   [van der Vaart and Wellner, 1996]

In the i.i.d. case with the true parametric model it holds

$$V_0^{-1} \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)(u_i - 1) \xrightarrow{\mathcal{L}^{\bullet}} \mathcal{N}(0, \boldsymbol{I}_p),$$

for almost every i.i.d. sequence $u_1, u_2 \ldots$ s.t. $\mathbb{E}^{\bullet} u_i = 1, \mathrm{Var}^{\bullet} u_i = 1$, with

$$V_0^2 \stackrel{\mathrm{def}}{=} \mathrm{Var}\{\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*)\}$$
$$= D_0^2 \text{ for the true parametric model.}$$

Therefore, in the i.i.d. parametric case the approximating vectors $\boldsymbol{\xi}$ and $\boldsymbol{\xi}^{\bullet} \approx \boldsymbol{\xi}^{\bullet}(\theta^*)$ have the same limit distributions.

Introduce for $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \boldsymbol{I}_p)$, fixed $\Gamma_0 = C\sqrt{p}$ and arbitrary $\boldsymbol{\gamma} \in \mathbb{R}^p$, $\|\boldsymbol{\gamma}\| = 1$:

$$h(\mu, t) \stackrel{\text{def}}{=} \exp(\mu t/2) \, \mathbb{P}\left(\|\boldsymbol{\varepsilon} + \sqrt{\mu t}\boldsymbol{\gamma}\| \le \mu^{-1/2}\Gamma_0\right).$$

---

**Theorem**

*It holds with probability* $\ge 1 - Ce^{-\mathbf{x}}$

$$\sup_{\mu \in (0,1)} \left| \frac{\mathbb{E}^{\bullet} h(\mu, \|\boldsymbol{\xi}^{\bullet}\|^2)}{\mathbb{E} h(\mu, \|\boldsymbol{\xi}\|^2)} - 1 \right| \le C\sqrt{\frac{(p + \mathbf{x})^3}{n}}.$$

---

Get to the linear exponent w.r.t. $\boldsymbol{\xi}$ by

$$\exp\left(\mu\|\boldsymbol{\xi}\|^2/2\right) I\!\!P\left(\|\boldsymbol{\varepsilon} + \mu^{1/2}\boldsymbol{\xi}\| \le \mu^{-1/2}\Gamma_0 \,|\, \boldsymbol{\xi}\right)$$

$$= \frac{1}{(2\pi\mu)^{p/2}} \int_{\|\boldsymbol{\gamma}\| \le \Gamma_0} \exp\left(\boldsymbol{\gamma}^\top\boldsymbol{\xi} - \frac{1}{2\mu}\|\boldsymbol{\gamma}\|^2\right) d\boldsymbol{\gamma}.$$

Use the Taylor expansion of $\log I\!\!E \exp\left(\lambda\boldsymbol{\gamma}^\top\boldsymbol{\xi}\right)$ w.r.t. $|\lambda| \le \Gamma_0 = \mathrm{C}\sqrt{p}$.

## A cumulative bound

Let $\mathfrak{z}_\alpha^{\bullet}$ denote the upper $\alpha$-quantile of $L^{\bullet}(\widetilde{\boldsymbol{\theta}}^{\bullet}, \widetilde{\boldsymbol{\theta}})$.

**Theorem**

*It holds with probability* $\geq 1 - C\mathrm{e}^{-\mathtt{y}}$

$$\mathbb{P}\left(L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) > \mathfrak{z}_\alpha^{\bullet} + \Delta_{cum}\right) - \alpha \leq \alpha\delta_F,$$

$$\mathbb{P}\left(L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) > \mathfrak{z}_\alpha^{\bullet} - \Delta_{cum}\right) - \alpha \geq -\alpha\delta_F,$$

*where*

$$\Delta_{cum}, \delta_F \lesssim \sqrt{\frac{(p + \mathtt{y})^3}{n}}.$$

Compare the approximating terms $I\!E\|\boldsymbol{\xi}\|^2$ and $I\!E^{\circ}\|\boldsymbol{\xi}^{\circ}\|^2$ :

$$I\!E\|\boldsymbol{\xi}\|^2 = \operatorname{tr}\left(D_0^{-1}V_0^2 D_0^{-1}\right), \qquad I\!E^{\circ}\|\boldsymbol{\xi}^{\circ}\|^2 = \operatorname{tr}\left(\mathfrak{D}_0^{-1}\mathcal{V}_0^2 \mathfrak{D}_0^{-1}\right).$$

$$
\begin{aligned}
V_0^2 &\stackrel{\text{def}}{=} \operatorname{Var}\nabla_{\boldsymbol{\theta}}L(\boldsymbol{\theta}^*) \\
&= \sum_{i=1}^{n} I\!E\left[\nabla_{\boldsymbol{\theta}}\ell_i(\boldsymbol{\theta}^*)\nabla_{\boldsymbol{\theta}}\ell_i(\boldsymbol{\theta}^*)^{\top}\right] - \sum_{i=1}^{n} I\!E\left[\nabla_{\boldsymbol{\theta}}\ell_i(\boldsymbol{\theta}^*)\right] I\!E\left[\nabla_{\boldsymbol{\theta}}\ell_i(\boldsymbol{\theta}^*)\right]^{\top}, \\
\mathcal{V}_0^2 &\stackrel{\text{def}}{=} \operatorname{Var}^{\circ}\nabla_{\boldsymbol{\theta}}L^{\circ}(\widetilde{\boldsymbol{\theta}}) \\
&= \sum_{i=1}^{n} \nabla_{\boldsymbol{\theta}}\ell_i(\widetilde{\boldsymbol{\theta}})\nabla_{\boldsymbol{\theta}}\ell_i(\widetilde{\boldsymbol{\theta}})^{\top}.
\end{aligned}
$$

The relation of the blue matrices in spectral norm is $\leq C\sqrt{(p+x)^3/n}$ . The magenta matrix adds the modelling bias, bounded by condition $(SmB)$ .

$(ED_3)$ *It holds for all* $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$, $\mathbf{r} \leq \mathbf{r}_0$ *and for* $j = 1, 2, 3$ *and* $|\lambda| \leq \mathbf{g}$

$$\sup_{\substack{\boldsymbol{\gamma}_j \in \mathbb{R}^p, \\ \|\boldsymbol{\gamma}_j\| \leq 1}} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega_1} \boldsymbol{\gamma}_3^\top \nabla_{\boldsymbol{\theta}} \left[ \boldsymbol{\gamma}_1^\top D_0^{-1} \nabla_{\boldsymbol{\theta}}^2 \zeta(\boldsymbol{\theta}) D_0^{-1} \boldsymbol{\gamma}_2 \right] \right\} \leq \frac{\nu_0^2 \lambda^2}{2}.$$

$(SmB)$ *There exists a constant* $\delta_\xi^2 \lesssim \sqrt{p/n^3}$ *such that it holds for all*
*$i = 1, \ldots, n$*

$$\left\| D_0^{-1} \mathbb{E} \nabla_{\boldsymbol{\theta}} \log \frac{dP_i(\boldsymbol{\theta}^*)}{d\mu_0}(Y_i) \right\| \leq \delta_\xi$$

$(SD_0)$ *There exists a constant $\delta_v \geq 0$ such that it holds for all $i = 1, \ldots, n$ with dominating probability*

$$\left\| H_0^{-1} \left\{ \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top - I\!E \left[ \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top \right] \right\} H_0^{-1} \right\| \leq \delta_v,$$

where

$$H_0^2 \stackrel{\text{def}}{=} \sum_{i=1}^{n} I\!E \left\{ \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*) \nabla_{\boldsymbol{\theta}} \ell_i(\boldsymbol{\theta}^*)^\top \right\}.$$

(Condition for the non-commutative Bernstein inequality by [Koltchinskii et al., 2011])

Bickel, P. J. and Kleijn, B. J. K. (2012).

The semiparametric Bernstein-von Mises theorem.

*Ann. Statist.*, 40(1):206–237.

Bontemps, D. (2011).

Bernstein-von Mises theorems for Gaussian regression with increasing number of regressors.

*Ann. Statist.*, 39(5):2557–2584.

Boucheron, S. and Gassiat, E. (2009).

A Bernstein-von Mises theorem for discrete probability distributions.

*Electron. J. Stat.*, 3:114–148.

Bunke, O. and Milhaud, X. (1998).

Asymptotic behavior of Bayes estimates under possibly incorrect models.

*Ann. Statist.*, 26(2):617–644.

Castillo, I. (2012).

A semiparametric Bernstein - von Mises theorem for Gaussian process priors.

*Probability Theory and Related Fields*, 152:53–99.

10.1007/s00440-010-0316-5.

Castillo, I. and Nickl, R. (2013).

Nonparametric Bernstein–von Mises theorems in Gaussian white noise.

*Ann. Statist.*, 41(4):1999–2028.

Castillo, I. and Rousseau, J. (2013).

A general bernstein–von mises theorem in semiparametric models.

Cox, D. D. (1993).

An analysis of Bayesian inference for nonparametric regression.

*Ann. Stat.*, 21(2):903–923.

Freedman, D. (1999).

On the Bernstein-von Mises theorem with infinite-dimensional parameters.

*Ann. Stat.*, 27(4):1119–1140.

Ghosal, S. (1999).

Asymptotic normality of posterior distributions in high-dimensional linear models.

*Bernoulli*, 5(2):315–331.

Ghosal, S. (2000).

Asymptotic normality of posterior distributions for exponential families when the number of parameters tends to infinity.

*J. Multivariate Anal.*, 74(1):49–68.

Ghosal, S., Ghosh, J. K., and van der Vaart, A. W. (2000).
Convergence rates of posterior distributions.
*Ann. Statist.*, 28(2):500–531.

Ghosal, S. and van der Vaart, A. (2007).
Convergence rates of posterior distributions for noniid observations.
*Ann. Statist.*, 35:192.

Kim, Y. (2006).
The Bernstein-von Mises theorem for the proportional hazard model.
*Ann. Statist.*, 34(4):1678–1700.

Kleijn, B. J. K. and van der Vaart, A. W. (2006).
Misspecification in infinite-dimensional Bayesian statistics.
*Ann. Statist.*, 34(2):837–877.

Kleijn, B. J. K. and van der Vaart, A. W. (2012).
The Bernstein-von-Mises theorem under misspecification.
*Electronic J. Statist.*, 6:354–381.

Koltchinskii, V., Lounici, K., and Tsybakov, A. B. (2011).

Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion.
*Ann. Stat.*, 39(5):2302–2329.

Rivoirard, V. and Rousseau, J. (2012).
Bernstein–von mises theorem for linear functionals of the density.
*Ann. Stat.*, 40(3):1489–1523.

Shen, X. (2002).
Asymptotic normality of semiparametric and nonparametric posterior distributions.
*J. Amer. Statist. Assoc.*, 97(457):222–235.

van der Vaart, A. and Wellner, J. A. (1996).
*Weak convergence and empirical processes. With applications to statistics.*
Springer Series in Statistics. New York, Springer.