**Weierstraß-Institut für**
**Angewandte Analysis und Stochastik**

# Fisher and Wilks expansions with applications to statistical inference

Vladimir Spokoiny,
WIAS, HU Berlin

# Outline

Data $\boldsymbol{Y} \sim I\!\!P$. Aim: infer on $I\!\!P$.

Parametric assumption (PA): $I\!\!P \in (I\!\!P_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta \subseteq I\!\!R^p) \ll \boldsymbol{\mu}_0$.

Maximum likelihood estimator (MLE):

$$\widetilde{\boldsymbol{\theta}} \stackrel{\mathrm{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \log \frac{d I\!\!P_{\boldsymbol{\theta}}}{d \boldsymbol{\mu}_0}(\boldsymbol{Y})$$

PA-PW: $I\!\!P \not\in (I\!\!P_{\boldsymbol{\theta}})$. Target of estimation ?

$$\boldsymbol{\theta}^* \stackrel{\mathrm{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} I\!\!E L(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} I\!\!E \log \frac{d I\!\!P_{\boldsymbol{\theta}}}{d I\!\!P} = \operatorname*{argmin}_{\boldsymbol{\theta} \in \Theta} \mathcal{K}(I\!\!P, I\!\!P_{\boldsymbol{\theta}}).$$

Under PA: $I\!\!P = I\!\!P_{\boldsymbol{\theta}^*}$ and

$$\operatorname*{argmin}_{\boldsymbol{\theta} \in \Theta} \mathcal{K}(I\!\!P_{\boldsymbol{\theta}^*}, I\!\!P_{\boldsymbol{\theta}}) = \boldsymbol{\theta}^*$$

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

---

**Theorem**

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

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

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

*with*

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

*Here* $\diamondsuit(\mathbf{x})$ *and* $\Delta(\mathbf{x})$ *are explicit error terms.*

---

Given

– $\boldsymbol{Y}$ , response,

– $\Sigma = \mathrm{Cov}(\boldsymbol{Y})$ , its covariance matrix

– $\Psi$ , design matrix of regressors:

$$\boldsymbol{Y} = \Psi^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon}, \qquad I\!\!E\boldsymbol{\varepsilon} = 0, \qquad \mathrm{Cov}(\boldsymbol{\varepsilon}) = \Sigma.$$

PA: $\boldsymbol{Y} \sim \mathcal{N}(\Psi^\top \boldsymbol{\theta}, \Sigma)$ :

$$L(\boldsymbol{\theta}) = -\frac{1}{2}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta}) + R$$

Study under true: $I\!\!E\boldsymbol{Y} = \boldsymbol{f}$ and $\mathrm{Cov}(\boldsymbol{Y}) = \Sigma_0$ .

$$L(\boldsymbol{\theta}) = -\frac{1}{2}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta}) + R,$$

### Lemma

$L(\boldsymbol{\theta})$ *is quadratic in* $\boldsymbol{\theta}$ *and it holds with* $\mathbb{E}\boldsymbol{Y} = \boldsymbol{f}$, $\boldsymbol{\varepsilon} \stackrel{\text{def}}{=} \boldsymbol{Y} - \boldsymbol{f}$:

$$\nabla^2 L(\boldsymbol{\theta}^*) = -\Psi \Sigma^{-1} \Psi^\top$$

$$D_0^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{\theta}^*) = \Psi \Sigma^{-1} \Psi^\top,$$

$$\widetilde{\boldsymbol{\theta}} = D_0^{-2} \Psi \Sigma^{-1} \boldsymbol{Y},$$

$$\boldsymbol{\theta}^* = D_0^{-2} \Psi \Sigma^{-1} \boldsymbol{f},$$

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1} \nabla L(\boldsymbol{\theta}^*) = D_0^{-1} \Psi \Sigma^{-1} \boldsymbol{\varepsilon},$$

$$D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \equiv \boldsymbol{\xi},$$

$$L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) \equiv \|\boldsymbol{\xi}\|^2/2.$$

$$L(\boldsymbol{\theta}) = -\frac{1}{2}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta})^\top \Sigma^{-1}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta}) + R,$$

$$\nabla L(\boldsymbol{\theta}) = \Psi \Sigma^{-1}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta}),$$

$$\nabla^2 L(\boldsymbol{\theta}) \equiv -\Psi \Sigma^{-1}\Psi^\top$$

$L(\boldsymbol{\theta})$ is quadratic in $\boldsymbol{\theta}$ and it holds with $I\!\!E \boldsymbol{Y} = \boldsymbol{f}$, $\boldsymbol{\varepsilon} \stackrel{\text{def}}{=} \boldsymbol{Y} - \boldsymbol{f}$:

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

$$\widetilde{\boldsymbol{\theta}} = D_0^{-2}\Psi \Sigma^{-1}\boldsymbol{Y}, \qquad \nabla L(\widetilde{\boldsymbol{\theta}}) = \Psi \Sigma^{-1}(\boldsymbol{Y} - \Psi^\top \widetilde{\boldsymbol{\theta}}) = 0,$$

$$\boldsymbol{\theta}^* = D_0^{-2}\Psi \Sigma^{-1}\boldsymbol{f}, \quad \nabla I\!\!E L(\boldsymbol{\theta}^*) = \Psi \Sigma^{-1}(\boldsymbol{f} - \Psi^\top \boldsymbol{\theta}^*) = 0,$$

$$\boldsymbol{\xi} \stackrel{\text{def}}{=} D_0^{-1}\nabla L(\boldsymbol{\theta}^*) \qquad\qquad = D_0^{-1}\Psi \Sigma^{-1}(\boldsymbol{Y} - \Psi^\top \boldsymbol{\theta}^*) = D_0^{-1}\Psi \Sigma^{-1}\boldsymbol{\varepsilon}.$$

Hence $D_0\big(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\big) = \boldsymbol{\xi}$ and

$$L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}^*) = -\frac{1}{2}(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top \nabla^2 L(\widetilde{\boldsymbol{\theta}})(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) = -\frac{1}{2}\big\|D_0\big(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\big)\big\|^2 = -\frac{1}{2}\|\boldsymbol{\xi}\|^2,$$

Under PA: $\boldsymbol{Y} \sim \mathcal{N}(\Psi^{\top}\boldsymbol{\theta}^{*}, \Sigma)$.

Then $\boldsymbol{\xi} = D_0^{-1}\Psi\Sigma^{-1}(\boldsymbol{Y} - I\!E\boldsymbol{Y})$ is normal zero mean and

$$\mathrm{Var}(\boldsymbol{\xi}) = \mathrm{Var}\big(D_0^{-1}\Psi\Sigma^{-1}\boldsymbol{\varepsilon}\big) = D_0^{-1}\Psi\Sigma^{-1}\,\mathrm{Var}(\boldsymbol{\varepsilon})\Sigma^{-1}\Psi D_0^{-1} = \boldsymbol{I}_p\,.$$

Therefore, $D_0\big(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^{*}\big) = \boldsymbol{\xi}$ is standard normal and

$$2L(\widetilde{\boldsymbol{\theta}}) - 2L(\boldsymbol{\theta}^{*}) = \|\boldsymbol{\xi}\|^2 \sim \chi_p^2$$

If $z_\alpha^2$ is the $1 - \alpha$ quantile of $\chi_p^2$, then

$$\mathcal{E}(z_\alpha) = \big\{\boldsymbol{\theta} \colon \|D_0\big(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}\big)\| \leq z_\alpha\big\} = \big\{\boldsymbol{\theta} \colon L(\widetilde{\boldsymbol{\theta}}) - L(\boldsymbol{\theta}) \leq z_\alpha^2/2\big\}$$

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

$$I\!P\big(\boldsymbol{\theta}^{*} \notin \mathcal{E}(z_\alpha)\big) = \alpha.$$

► The Fisher and Wilks expansions are only based on geometric features of the likelihood ( $L(\boldsymbol{\theta})$ is quadratic in $\boldsymbol{\theta}$ ).

► The true distribution is not involved.

► Applies for any sample size.

► For inference, the PA is important. It only concerns the distribution of $\boldsymbol{\xi}$ .

► PA-PW: Let $\mathrm{Var}(\boldsymbol{Y}) = \Sigma_0 \neq \Sigma$ . Then with $D_0^2 = \Psi \Sigma^{-1} \Psi^\top$

$$\mathrm{Var}\big\{\nabla L(\boldsymbol{\theta}^*)\big\} = \mathrm{Var}\big\{\Psi \Sigma^{-1} \boldsymbol{Y}\big\} = \Psi \Sigma^{-1} \Sigma_0 \Sigma^{-1} \Psi^\top \stackrel{\text{def}}{=} V_0^2 \neq D_0^2$$

and (the sandwich formula)

$$\mathrm{Var}(\boldsymbol{\xi}) = \mathrm{Var}\big\{D_0^{-1} \nabla L(\boldsymbol{\theta}^*)\big\} = D_0^{-1} V_0^2 D_0^{-1} \neq \boldsymbol{I}_p.$$

Ley $\boldsymbol{Y} = (Y_1, \ldots, Y_n)$ be i.i.d. from $P$.

PA: $P \in (P_{\boldsymbol{\theta}}, \boldsymbol{\theta} \in \Theta)$, a regular family with $\ell(y, \boldsymbol{\theta}) = \log p(y, \boldsymbol{\theta})$.

$$L(\boldsymbol{\theta}) = \sum_{i=1}^{n} \ell(Y_i, \boldsymbol{\theta}), \qquad \widetilde{\boldsymbol{\theta}}_n = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \, L(\boldsymbol{\theta}).$$

## Theorem

*Assume* PA: $P = P_{\boldsymbol{\theta}^*} \in (P_{\boldsymbol{\theta}})$. *Then*

$$\sqrt{n \mathbb{F}_{\boldsymbol{\theta}^*}} \left( \widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^* \right) \xrightarrow{w} \mathcal{N}(0, I_p),$$

$$L(\widetilde{\boldsymbol{\theta}}_n) - L(\boldsymbol{\theta}^*) \xrightarrow{w} \chi_p^2 / 2$$

*where* $\mathbb{F}_{\boldsymbol{\theta}^*}$ *is the Fisher information matrix:*

$$\mathbb{F}_{\boldsymbol{\theta}^*} = -\nabla^2 E \, \ell(Y_1, \boldsymbol{\theta}^*) = \operatorname{Var}\{\nabla \ell(Y_1, \boldsymbol{\theta}^*)\}.$$

(Non-asymptotic) expansions:

$$\left\| D_0(\widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_n \right\| \leq \diamondsuit(\mathbf{x}),$$

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

where

$$D_0^2 \;=\; D_n^2 = -n\nabla^2 E\,\ell(Y_1, \boldsymbol{\theta}^*) = n\mathbb{F}_{\boldsymbol{\theta}^*}$$

$$\boldsymbol{\xi} \;=\; \boldsymbol{\xi}_n = (n\mathbb{F}_{\boldsymbol{\theta}^*})^{-1/2} \sum_{i=1}^{n} \nabla\ell(Y_i, \boldsymbol{\theta}^*)$$

Under PA $\nabla\ell(Y_i, \boldsymbol{\theta}^*)$ are i.i.d. zero mean with $\mathrm{Var}\big\{\nabla\ell(Y_1, \boldsymbol{\theta}^*)\big\} = \mathbb{F}_{\boldsymbol{\theta}^*}$, and by CLT

$$\boldsymbol{\xi}_n \xrightarrow{w} \mathcal{N}(0, \boldsymbol{I}_p)$$

For

$$\widetilde{\boldsymbol{\theta}}_n = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^{n} \ell(Y_i, \boldsymbol{\theta}),$$

it holds with $D_n^2 = n \mathbb{F}_{\boldsymbol{\theta}^*}$

$$\left\| D_n(\widetilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_n \right\| \leq \diamondsuit_n(\mathbf{x}),$$

$$\left| L(\widetilde{\boldsymbol{\theta}}_n) - L(\boldsymbol{\theta}^*) - \frac{\|\boldsymbol{\xi}_n\|^2}{2} \right| \leq \Delta_n(\mathbf{x}).$$

The error terms satisfy

$$\diamondsuit_n(\mathbf{x}) \leq \mathtt{C} \sqrt{\frac{(p + \mathbf{x})^2}{n}}, \qquad \Delta_n(\mathbf{x}) \leq \mathtt{C} \sqrt{\frac{(p + \mathbf{x})^3}{n}}.$$

and

$$\|\boldsymbol{\xi}_n\|^2 \leq p + \mathtt{C}\mathbf{x}.$$

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:

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

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

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

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

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

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

[Portnoy, 1984]: M-estimator i.i.d. or linear models:

– $p_n \log(p_n)/n \to 0$, consistency;

– $p_n^2 \log^2(p)/n \to 0$, asymptotic normality; (a counterexample for $p^2/n \to \infty$).

[Portnoy, 1988]: MLE for a GLM:

– $p_n^{3/2} \log(n)/n \to 0$, Wilks Theorem $p_n^{-1/2} L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - p_n^{1/2} \xrightarrow{w} \mathcal{N}(0,1)$;

Sieve estimation:

[Birgé and Massart, 1993], [Chen, 1993, 1997], [Van de Geer, 1993, van de Geer, 2002]; . . .

# Outline

Aim:

- minimal non-restrictive and natural conditions

- possibly sharp bounds

- all constants explicit, no asymptotic arguments

- model misspecification incorporated

- self-contained

- Concentration and large deviations: for some $\mathbf{r}_0$

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

  where $\quad \Theta_0(\mathbf{r}) \stackrel{\text{def}}{=} \big\{ \boldsymbol{\theta} \colon \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r} \big\}.$

- Local quadratic approximation of the expected log-likelihood:

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \frac{2\mathbb{E}L(\boldsymbol{\theta}^*) - 2\mathbb{E}L(\boldsymbol{\theta})}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2} \leq \delta(\mathbf{r}).$$

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

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \big| D_0^{-1}\big\{ \nabla\zeta(\boldsymbol{\theta}) - \nabla\zeta(\boldsymbol{\theta}^*) \big\} \big| \leq \varrho(\mathbf{r}, \mathbf{x}).$$

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

## Local quadraticity of $I\!E L(\boldsymbol{\theta})$

Define

$$D^2(\boldsymbol{\theta}) \stackrel{\mathrm{def}}{=} -\nabla^2 I\!E L(\boldsymbol{\theta}).$$

Then $D_0^2 = D_0^2(\boldsymbol{\theta}^*)$.

$(\mathcal{L}_0)$  *For each $\mathbf{r} \leq \mathbf{r}_0$, there is a constant $\delta(\mathbf{r}) \leq 1/2$ such that it holds for any*
$\boldsymbol{\theta} \in \Theta_0(\mathbf{r}) = \{\boldsymbol{\theta} \in \Theta \colon \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}\}$ :

$$\left\| D_0^{-1} D_0^2(\boldsymbol{\theta}) D_0^{-1} - I_p \right\|_\infty \leq \delta(\mathbf{r}).$$

By the second order Taylor expansion at $\boldsymbol{\theta}^*$ for any $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ :

$$\left| -2I\!E L(\boldsymbol{\theta}) + 2I\!E L(\boldsymbol{\theta}^*) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|^2 \right| \leq \delta(\mathbf{r})\mathbf{r}^2,$$

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

$$\leq \left\| \{I_p - D_0^{-1} D^2(\boldsymbol{\theta}^\circ) D_0^{-1}\} D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| \leq \delta(\mathbf{r})\mathbf{r}.$$

Aim: To bound the error of the local constant approximation of the gradient (vector) process

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\| D_0^{-1} \left\{ \nabla\zeta(\boldsymbol{\theta}) - \nabla\zeta(\boldsymbol{\theta}^*) \right\} \right\|$$

$(ED_2)$  *There exist a value $\omega > 0$ and for each $\mathbf{r} > 0$, a constant $\mathbf{g}(\mathbf{r}) > 0$ such that*
$\zeta(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - I\!E L(\boldsymbol{\theta})$ *satisfies for any $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$ :*

$$\sup_{\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in I\!R^p} \log I\!E \exp\left\{ \frac{\lambda}{\omega} \, \frac{\boldsymbol{\gamma}_1^\top \nabla^2 \zeta(\boldsymbol{\theta}) \boldsymbol{\gamma}_2}{\|D_0 \boldsymbol{\gamma}_1\| \cdot \|D_0 \boldsymbol{\gamma}_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq \mathbf{g}(\mathbf{r}).$$

Meaning: The second derivative of $\zeta(\boldsymbol{\theta})$ w.r.t. the local argument $\boldsymbol{\upsilon} = D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ is small.

Usually $\omega \asymp \|D_0^{-1}\| \asymp n^{-1/2}$.

Use $\boldsymbol{v} = D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ and consider $\mathcal{Y}(\boldsymbol{v}) = \omega^{-1}D_0^{-1}\{\nabla\zeta(\boldsymbol{\theta}) - \nabla\zeta(\boldsymbol{\theta}^*)\}$ :

$$\sup_{\boldsymbol{\theta}\in\Theta_0(\mathbf{r})}\big\|D_0^{-1}\{\nabla\zeta(\boldsymbol{\theta}) - \nabla\zeta(\boldsymbol{\theta}^*)\}\big\| \;=\; \omega\sup_{\boldsymbol{v}\in\Upsilon_\circ(\mathbf{r})}\|\mathcal{Y}(\boldsymbol{v})\|,$$

$$\Upsilon_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{v}\colon \|\boldsymbol{v}\| \leq \mathbf{r}\}.$$

For any $\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2 \in I\!\!R^p$ with $\|\boldsymbol{\gamma}_1\| = \|\boldsymbol{\gamma}_2\| = 1$, condition $(ED_2)$ implies

$$\log I\!\!E \exp\left\{\lambda\boldsymbol{\gamma}_1^\top\nabla\mathcal{Y}(\boldsymbol{v})\boldsymbol{\gamma}_2\right\} \;=\; \log I\!\!E \exp\left\{\frac{\lambda}{\omega}\boldsymbol{\gamma}_1^\top D_0^{-1}\nabla^2\zeta(\boldsymbol{\theta})D_0^{-1}\boldsymbol{\gamma}_2\right\} \leq \frac{\nu_0^2\lambda^2}{2}.$$

## A bound for the norm of a vector stochastic process

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

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

---

### 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:*

$$z_{\mathbb{H}}(\mathbf{x}) = \begin{cases} \sqrt{\mathbb{H}_2 + 2\mathbf{x}}, & \text{if } \mathbb{H}_2 + 2\mathbf{x} \le \mathbf{g}^2, \\ \mathbf{g}^{-1}\mathbf{x} + \frac{1}{2}\big(\mathbf{g}^{-1}\mathbb{H}_2 + \mathbf{g}\big), & \text{if } \mathbb{H}_2 + 2\mathbf{x} > \mathbf{g}^2. \end{cases}$$

*Here* $\mathbb{H}_2 = 4p$ *and* $\mathbb{H}_1 = 2p^{1/2}$, $\mathbf{g} = \mathbf{g}(\mathbf{r})$.

**Local linear approximation of the gradient**

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

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

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

---

**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 $z_{\mathbb{H}}(\mathbf{x})$.

Define

$$\chi(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \stackrel{\text{def}}{=} D_0^{-1} \{ \nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^*) + D_0^2 (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \}.$$

By Theorem 5

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} \big\| \chi(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \big\| \leq \Diamond(\mathbf{r}_0, \mathbf{x}).$$

Suppose that $\widetilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}_0)$ on $\Omega(\mathbf{x})$. Then

$$\big\| D_0^{-1} \{ \nabla L(\widetilde{\boldsymbol{\theta}}) - \nabla L(\boldsymbol{\theta}^*) \} + D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \big\| \leq \Diamond(\mathbf{r}, \mathbf{x}).$$

The use of $\nabla L(\widetilde{\boldsymbol{\theta}}) = 0$ yields the Fisher expansion.

**A quadratic approximation**

Define $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2/2$ and

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) + \frac{1}{2}\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2$$

$$= L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ), \qquad \boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$$

With $\boldsymbol{\theta}^\circ$ fixed, the gradient $\nabla \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} \frac{d}{d\boldsymbol{\theta}}\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ)$ fulfills

$$\nabla \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = \nabla L(\boldsymbol{\theta}) - \nabla L(\boldsymbol{\theta}^\circ) + D_0^2(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ) = D_0\, \chi(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ);$$

This implies

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla \alpha(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ),$$

where $\boldsymbol{\theta}'$ is a point on the line connecting $\boldsymbol{\theta}$ and $\boldsymbol{\theta}^\circ$ and

$$\left| \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \right| = \left| (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top D_0 D_0^{-1} \nabla \alpha(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ) \right| \leq \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\| \sup_{\boldsymbol{\theta}' \in \Theta_0(\mathbf{r})} \left| \chi(\boldsymbol{\theta}', \boldsymbol{\theta}^\circ) \right|.$$

$$\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) + \frac{1}{2}\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2$$

$$= L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ), \qquad \boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$$

**Theorem**

*Suppose* $(\mathcal{L}_0)$ *,* $(ED_0)$ *, and* $(ED_2)$ *. For each* $\mathbf{r}$ *, it holds on a random set* $\Omega(\mathbf{r}, \mathbf{x})$ *of a dominating probability at least* $1 - \mathrm{e}^{-\mathbf{x}}$ *, it holds with any* $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$

$$\frac{|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} \le \Diamond(\mathbf{r}, \mathbf{x}), \qquad |\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \le \mathbf{r}\Diamond(\mathbf{r}, \mathbf{x}),$$

$$\frac{|\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta})|}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} \le \Diamond(\mathbf{r}, \mathbf{x}), \qquad |\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta})| \le \mathbf{r}\Diamond(\mathbf{r}, \mathbf{x}).$$

Let $\widetilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}_0)$ on $\Omega(\mathbf{x})$. For $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) = (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) - \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2/2$

$$\left| \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \right| = \left| L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \right| \leq \mathbf{r}_0 \Diamond(\mathbf{r}_0, \mathbf{x}), \qquad \boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r}_0)$$

The special case with $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ and $\boldsymbol{\theta}^\circ = \widetilde{\boldsymbol{\theta}}$ yields in view of $\nabla L(\widetilde{\boldsymbol{\theta}}) = 0$ for $\mathbf{r} = \mathbf{r}_0$

$$\left| L(\boldsymbol{\theta}^*) - L(\widetilde{\boldsymbol{\theta}}) + \|D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2/2 \right| = \left| \alpha(\boldsymbol{\theta}^*, \widetilde{\boldsymbol{\theta}}) \right| \leq \mathbf{r}_0 \Diamond(\mathbf{r}_0, \mathbf{x}). \tag{1}$$

Further, on the set of a dominating probability, it holds $\|\boldsymbol{\xi}\| \leq z(B, \mathbf{x})$ (later). Now

$$\begin{aligned}
& \left| \|D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2 - \|\boldsymbol{\xi}\|^2 \right| \\
& \leq 2 \|\boldsymbol{\xi}\| \cdot \|D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| + \|D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\|^2 \\
& \leq 2 z(B, \mathbf{x}) \Diamond(\mathbf{r}_0, \mathbf{x}) + \Diamond^2(\mathbf{r}_0, \mathbf{x}).
\end{aligned}$$

Together with (1), this yields

$$\left| L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|\boldsymbol{\xi}\|^2/2 \right| \leq \left\{ \mathbf{r}_0 + z(B, \mathbf{x}) \right\} \Diamond(\mathbf{r}_0, \mathbf{x}) + \Diamond^2(\mathbf{r}_0, \mathbf{x})/2.$$

The error term can be improved if the squared root of the excess is considered.

Indeed, if $\widetilde{\boldsymbol{\theta}} \in \Theta_0(\mathbf{r}_0)$

$$\left| \left\{ 2L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \right\}^{1/2} - \| D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \| \right| \leq \frac{\left| 2L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \| D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \|^2 \right|}{\| D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \|}$$

$$\leq \frac{2 \left| \alpha(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \right|}{\| D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \|} \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} \frac{2 \left| \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \right|}{\| D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \|} \leq 2 \diamondsuit(\mathbf{r}_0, \mathbf{x}).$$

The Fisher expansion allows to replace here the norm of the standardized error $D_0(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ with the norm of the normalized score $\boldsymbol{\xi}$.

Aim: find $\mathbf{r}_0$ ensuring

$$IP\big(\widetilde{\boldsymbol{\theta}} \notin \Theta_0(\mathbf{r}_0)\big) \leq C e^{-\mathbf{x}}.$$

▶ By definition $\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq 0$ . Suffices to check that

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) < 0 \qquad \forall \boldsymbol{\theta} \in \Theta \setminus \Theta_0(\mathbf{r}_0)$$

▶ Use the decomposition

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = IEL(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) + \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*)$$

▶ Bound $\|\boldsymbol{\xi}\| = \|D_0^{-1} \nabla \zeta(\boldsymbol{\theta}^*)\|$;
▶ Upper function device for the remainder

$$\sup_{\boldsymbol{\theta} \in \Theta \setminus \Theta_0(\mathbf{r}_0)} \big\{ \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) - f(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \big\} \leq 0 \qquad \text{w.h.p.}$$

$(\mathcal{L})$   *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*

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

**Theorem**

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

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

*where*

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

*Then*

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

The radius $\mathbf{r}_0$ has to fulfill

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

One can use that

► $\mathbf{b}(\mathbf{r}_0) \geq 1 - \delta(\mathbf{r}_0) \approx 1$,

► the constant $\omega$ and thus, $\varrho(\mathbf{r},\mathbf{x})$, is small, and

► $\mathbf{rb}(\mathbf{r})$ grows with $\mathbf{r}$.

A simple rule $\mathbf{r}_0 \geq (2+\delta)z(B,\mathbf{x})$ for some $\delta > 0$ works in most of cases.

($ED_0$)  *There exist a positive symmetric matrix $V_0^2$, and constants $\mathrm{g} > 0$, $\nu_0 \geq 1$ such that*
$\mathrm{Var}\{\nabla\zeta(\boldsymbol{\theta}^*)\} \leq V_0^2$ *and*

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top V_0^{-1} \nabla\zeta(\boldsymbol{\theta}^*)) \leq \frac{\nu_0^2 \|\boldsymbol{\gamma}\|^2}{2}, \qquad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq \mathrm{g}.$$

With $\boldsymbol{\eta} = V_0^1 \nabla\zeta(\boldsymbol{\theta}^*)$, it holds $\boldsymbol{\xi} = D_0^{-1} V_0 \boldsymbol{\eta}$ and

$$\|\boldsymbol{\xi}\|^2 = \boldsymbol{\eta}^\top B \boldsymbol{\eta}$$

for $B = D_0^{-1} V_0^2 D_0^{-1}$. Also define

$$\mathrm{p}_B \stackrel{\mathrm{def}}{=} \mathrm{tr}(B), \qquad \mathrm{v}_B^2 \stackrel{\mathrm{def}}{=} 2\,\mathrm{tr}(B^2), \qquad \lambda_B \stackrel{\mathrm{def}}{=} \lambda_{\max}(B).$$

Note that $\mathrm{p}_B = \mathbb{E}\|\boldsymbol{\xi}\|^2$. Moreover, if $\boldsymbol{\xi}$ is a Gaussian vector then $\mathrm{v}_B^2 = \mathrm{Var}(\|\boldsymbol{\xi}\|^2)$. If $V_0^2 = D_0^2$, then $\lambda_B = 1$.

Define $\mu_c = 2/3$, $\mathtt{p}_B = \operatorname{tr}(B)$, $\mathtt{v}_B^2 = 2\operatorname{tr}(B^2)$, and $\lambda_B = \lambda_{\max}(B)$

$$\mathtt{g}_c \stackrel{\text{def}}{=} \sqrt{\mathtt{g}^2 - \mu_c \mathtt{p}_B},$$

$$2\mathtt{x}_c \stackrel{\text{def}}{=} (\mathtt{g}^2/\mu_c - \mathtt{p}_B)/\lambda_B + \log\det\!\big(\boldsymbol{I}_p - \mu_c B/\lambda_B\big). \tag{3}$$

---

**Theorem (SP2012)**

Let $(ED_0)$ hold with $\nu_0 = 1$ and $\mathtt{g}^2 \geq 2\mathtt{p}_B$ . Then for each $\mathtt{x} > 0$

$$I\!P\big(\|\boldsymbol{\xi}\| \geq z(B,\mathtt{x})\big) = I\!P\big(\|B^{1/2}\boldsymbol{\eta}\| \geq z(B,\mathtt{x})\big) \;\leq\; 2\mathrm{e}^{-\mathtt{x}} + 8.4\mathrm{e}^{-\mathtt{x}_c},$$

where $z(B,\mathtt{x})$ is defined with $\mathtt{y}_c^2 \leq \mathtt{p}_B + 6\lambda_B \mathtt{x}_c$ by

$$z^2(B,\mathtt{x}) \stackrel{\text{def}}{=} \begin{cases} \mathtt{p}_B + 2\mathtt{v}_B \mathtt{x}^{1/2}, & \mathtt{x} \leq \mathtt{v}_B/(18\lambda_B), \\ \mathtt{p}_B + 6\lambda_B \mathtt{x}, & \mathtt{v}_B/(18\lambda_B) < \mathtt{x} \leq \mathtt{x}_c, \\ \big|\mathtt{y}_c + 2\lambda_B(\mathtt{x} - \mathtt{x}_c)/\mathtt{g}_c\big|^2, & \mathtt{x} > \mathtt{x}_c. \end{cases}$$

$$\mathtt{p}_B \;=\; \mathrm{tr}\big(B\big), \qquad \mathtt{v}_B^2 = 2\,\mathrm{tr}\big(B^2\big), \qquad \lambda_B = \lambda_{\max}\big(B\big).$$

$$z^2(B,\mathtt{x}) \;\overset{\mathrm{def}}{=}\; \begin{cases} \mathtt{p}_B + 2\mathtt{v}_B \mathtt{x}^{1/2}, & \mathtt{x} \le \mathtt{v}_B/(18\lambda_B), \\ \mathtt{p}_B + 6\lambda_B \mathtt{x}, & \mathtt{v}_B/(18\lambda_B) < \mathtt{x} \le \mathtt{x}_c, \\ \big| \mathtt{y}_c + 2\lambda_B(\mathtt{x}-\mathtt{x}_c)/\mathtt{g}_c \big|^2, & \mathtt{x} > \mathtt{x}_c. \end{cases}$$

Depending on the value $\mathtt{x}$ , we observe three types of tail behavior of the quadratic form $\|\boldsymbol{\xi}\|^2$ :

- The sub-Gaussian regime for $\mathtt{x} \le \mathtt{v}_B/(18\lambda_B)$

- The Poissonian regime for $\mathtt{x} \le \mathtt{x}_c$

- The value $\mathtt{x}_c$ from (3) is of order $\mathtt{g}^2$ . In all our results we suppose that $\mathtt{g}^2$ and hence, $\mathtt{x}_c$ is sufficiently large;

The quadratic form $\|\boldsymbol{\xi}\|^2$ can be bounded with a dominating probability by $\mathtt{p}_B + 6\lambda_B \mathtt{x}$ for a proper $\mathtt{x}$ .

## A "squared norm" trick

Let $\boldsymbol{\xi}$ be a random vector in $\mathbb{R}^p$ satisfying the condition

$$\log \mathbb{E} \exp(\boldsymbol{\gamma}^\top \boldsymbol{\xi}) \leq \frac{\nu_0^2 \|\boldsymbol{\gamma}\|^2}{2}, \qquad \boldsymbol{\gamma} \in \mathbb{R}^p, \|\boldsymbol{\gamma}\| \leq \mathsf{g}.$$

For simplicity we take here $B = 1$.

Aim: to bound $\|\boldsymbol{\xi}\|^2$.

A sup-representation:

$$\|\boldsymbol{\xi}\|^2 = \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p} \{\boldsymbol{\gamma}^\top \boldsymbol{\xi} - \|\boldsymbol{\gamma}\|^2/2\}, \qquad \|\boldsymbol{\xi}\| = \sup_{\boldsymbol{\gamma} \in \mathbb{R}^p : \|\boldsymbol{\gamma}\| \leq 1} \boldsymbol{\gamma}^\top \boldsymbol{\xi}.$$

Too rough to get a sharp bound on $\|\boldsymbol{\xi}\|$ with entropy arguments.

An exp-representation: for any $\mu < 1$

$$\exp\{\mu \|\boldsymbol{\xi}\|^2/2\} = \mathsf{C}_p(\mu) \int_{\mathbb{R}^p} \exp\{\boldsymbol{\gamma}^\top \boldsymbol{\xi} - \|\boldsymbol{\gamma}\|^2/(2\mu)\} d\boldsymbol{\gamma}$$

**An upper function for the stochastic component** $\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - I\!\!EL(\boldsymbol{\theta})$

The proof is based on the following bound: for each $\mathbf{r}$

$$I\!\!P\left(\sup_{\boldsymbol{\theta}\in\Theta_0(\mathbf{r})}\left|\zeta(\boldsymbol{\theta}) - \zeta(\boldsymbol{\theta}^*) - (\boldsymbol{\theta}-\boldsymbol{\theta}^*)^\top\nabla\zeta(\boldsymbol{\theta}^*)\right| \geq 3\nu_0\, z_{\mathbb{H}}(\mathbf{x})\,\omega\,\mathbf{r}\right) \leq \mathrm{e}^{-\mathbf{x}}.$$

This bound is a special case of the general result from Theorem 9 below. It implies by Theorem 10 with $\rho = 1/2$ on a set $\Omega(\mathbf{x})$ of probability at least $1 - \mathrm{e}^{-\mathbf{x}}$ that for all $\mathbf{r} \geq \mathbf{r}_0$ and all $\boldsymbol{\theta}$ with $\|D_0(\boldsymbol{\theta}-\boldsymbol{\theta}^*)\| \leq \mathbf{r}$

$$\left|\zeta(\boldsymbol{\theta},\boldsymbol{\theta}^*) - (\boldsymbol{\theta}-\boldsymbol{\theta}^*)^\top\nabla\zeta(\boldsymbol{\theta}^*)\right| \leq \varrho(\mathbf{r},\mathbf{x})\,\mathbf{r},$$

where

$$\varrho(\mathbf{r},\mathbf{x}) = 6\nu_0\, z_{\mathbb{H}}\big(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)\big)\,\omega\,. \tag{4}$$

The use of $\nabla I\!\!EL(\boldsymbol{\theta}^*) = 0$ yields

$$\sup_{\boldsymbol{\theta}\in\Theta_0(\mathbf{r})}\left|L(\boldsymbol{\theta},\boldsymbol{\theta}^*) - I\!\!EL(\boldsymbol{\theta},\boldsymbol{\theta}^*) - (\boldsymbol{\theta}-\boldsymbol{\theta}^*)^\top\nabla L(\boldsymbol{\theta}^*)\right| \leq \varrho(\mathbf{r},\mathbf{x})\,\mathbf{r}.$$

By definition $\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq 0$. So, it suffices to check that $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) < 0$ for all $\boldsymbol{\theta} \in \Theta \setminus \Theta_0(\mathbf{r}_0)$.

We know

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left| L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - I\!\!E L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) \right| \leq \varrho(\mathbf{r}, \mathbf{x}) \, \mathbf{r}.$$

Also $\|\boldsymbol{\xi}\| \leq z(B, \mathbf{x})$ on $\Omega(\mathbf{x})$ and for each $\mathbf{r} \geq \mathbf{r}_0$

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left| (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla L(\boldsymbol{\theta}^*) \right|$$
$$\leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \times \|D_0^{-1} \nabla L(\boldsymbol{\theta}^*)\| = \mathbf{r}\|\boldsymbol{\xi}\| \leq z(B, \mathbf{x}) \, \mathbf{r}.$$

Condition $(\mathcal{L})$ implies $-2 I\!\!E L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \geq \mathbf{r}^2 \mathbf{b}(\mathbf{r})$ for each $\boldsymbol{\theta}$ with $\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| = \mathbf{r}$. We conclude that the condition

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

ensure $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) < 0$ for all $\boldsymbol{\theta} \notin \Theta_0(\mathbf{r}_0)$ with a dominating probability.

Let $\mathcal{U}(\boldsymbol{v})$ be a smooth stochastic process on an open subset $\Upsilon \subseteq \mathbb{R}^p$, and $\mathbb{E}\mathcal{U}(\boldsymbol{v}) \equiv 0$.

$(\mathcal{E}D)$   There exist $\mathrm{g} > 0$, $\nu_0 \geq 1$, and a symmetric $H_0 \geq 0$ s.t. it holds

$$\sup_{\boldsymbol{\gamma} \in \mathbb{R}^p \,:\, \|\boldsymbol{\gamma}\|=1} \log \mathbb{E} \exp \left\{ \lambda \frac{\boldsymbol{\gamma}^\top \nabla \mathcal{U}(\boldsymbol{v})}{\|H_0 \boldsymbol{\gamma}\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \qquad |\lambda| \leq \mathrm{g}.$$

We consider the local sets of the elliptic form $\Upsilon_\circ(\mathbf{r}) \stackrel{\mathrm{def}}{=} \left\{ \boldsymbol{v} : \|H_0(\boldsymbol{v} - \boldsymbol{v}_0)\| \leq \mathbf{r} \right\}$.

**Theorem**

*Let $(\mathcal{E}D)$ hold with some $\mathrm{g} > 0$, and a matrix $H_0$. For any $\mathrm{x} \geq 0$ and any $\mathbf{r} > 0$*

$$\mathbb{P} \left\{ \sup_{\boldsymbol{v} \in \Upsilon_\circ(\mathbf{r})} \left| \mathcal{U}(\boldsymbol{v}) - \mathcal{U}(\boldsymbol{v}_0) \right| \geq 3\nu_0 \, \mathbf{r} \, z_{\mathbb{H}}(\mathrm{x}) \right\} \leq \mathrm{e}^{-\mathrm{x}},$$

*where $z_{\mathbb{H}}(\mathrm{x})$ is given by the following rule: with $\mathbb{H} = 4p$*

$$z_{\mathbb{H}}(\mathrm{x}) = \begin{cases} \sqrt{\mathbb{H} + 2\mathrm{x}} & \text{if } \mathbb{H} + 2\mathrm{x} \leq \mathrm{g}^2, \\ \mathrm{g}^{-1}\mathrm{x} + \frac{1}{2}\left(\mathrm{g}^{-1}\mathbb{H} + \mathrm{g}\right) & \text{if } \mathbb{H} + 2\mathrm{x} > \mathrm{g}^2, \end{cases}$$

**Tools. An "upper function" device**

On $\Omega(\mathbf{r}, \mathbf{x})$, one can bound $\mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) \stackrel{\text{def}}{=} \mathcal{U}(\boldsymbol{v}) - \mathcal{U}(\boldsymbol{v}_0)$:

$$\left| \mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) \right| \leq 3\nu_0 \, \mathbf{r} \, z_{\mathbb{H}}(\mathbf{x}).$$

Aim: to build an upper function $f(\cdot)$ s.t. $\mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) - f(\boldsymbol{v}, \boldsymbol{v}_0)$ is bounded uniformly in all $\boldsymbol{v}$.

---

**Theorem**

*Let $(\mathcal{E}D)$ hold on $\mathcal{B}_{\mathbf{r}^*}(\boldsymbol{v}_0)$. Given $\mathbf{r}_0 < \mathbf{r}^*$, define $f(\mathbf{r}, \mathbf{r}_0)$ for some $\rho < 1$ as*

$$f(\mathbf{r}, \mathbf{r}_0) = 3\nu_0 \mathbf{r} \, z_{\mathbb{H}}\big(\mathbf{x} + \log(\mathbf{r}/\mathbf{r}_0)\big), \quad \mathbf{r}_0 \leq \mathbf{r} \leq \mathbf{r}^*. \tag{5}$$

*Then it holds*

$$\mathbb{P}\left( \sup_{\mathbf{r}_0 \leq \mathbf{r} \leq \mathbf{r}^*} \sup_{\boldsymbol{v} \in \Upsilon_\circ(\mathbf{r})} \left\{ \mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) - f\big(\rho^{-1}\mathbf{r}, \mathbf{r}_0\big) \right\} \geq 0 \right) \leq \frac{\rho}{1-\rho} e^{-\mathbf{x}}.$$

---

If $\mathbf{g} = \infty$, then $z_{\mathbb{H}}(\mathbf{x}) = \sqrt{2\mathbf{x} + 4p}$ and ($\rho = 1/2$)

$$f(\mathbf{r}, \mathbf{r}_0) = 3\nu_0 \mathbf{r} \sqrt{2\mathbf{x} + 4p + 2\log(\mathbf{r}/\mathbf{r}_0)}.$$

Idea: split $\mathcal{B}_{r^*}(\boldsymbol{v}_0)$ into slices $\mathcal{B}_{r_k}(\boldsymbol{v}_0) \setminus \mathcal{B}_{r_{k-1}}(\boldsymbol{v}_0)$ and apply Theorem 9 to each slice. By (5) and Theorem 9 for any $r > r_0$

$$\mathbb{P}\bigg( \sup_{\boldsymbol{v} \in \mathcal{B}_r(\boldsymbol{v}_0) \setminus \mathcal{B}_{\rho r}(\boldsymbol{v}_0)} \big\{ \mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) - f(r, r_0) \big\} \geq 0 \bigg)$$

$$\leq \mathbb{P}\bigg( \frac{1}{3\nu_0 r} \sup_{\boldsymbol{v} \in \mathcal{B}_r(\boldsymbol{v}_0)} \mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) \geq z_{\mathbb{H}}\big( x + \log(r/r_0) \big) \bigg) \leq \frac{r_0}{r} e^{-x}. \tag{6}$$

Define $r_k = r_0 \rho^{-k}$ for $k = 0, 1, 2, \ldots$ and $k^* \stackrel{\text{def}}{=} \log(r^*/r_0) + 1$. By (6)

$$\mathbb{P}\bigg( \sup_{\boldsymbol{v} \in \mathcal{B}_{r^*}(\boldsymbol{v}_0) \setminus \mathcal{B}_{r_0}(\boldsymbol{v}_0)} \big\{ \mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) - f\big( \rho^{-1} d(\boldsymbol{v}, \boldsymbol{v}_0), r_0 \big) \big\} \geq 0 \bigg)$$

$$\leq \sum_{k=1}^{k^*} \mathbb{P}\bigg( \frac{1}{r_k} \sup_{\boldsymbol{v} \in \mathcal{B}_{r_k}(\boldsymbol{v}_0) \setminus \mathcal{B}_{r_{k-1}}(\boldsymbol{v}_0)} \big\{ \mathcal{U}(\boldsymbol{v}, \boldsymbol{v}_0) - f(r_k, r_0) \big\} \geq 0 \bigg)$$

$$\leq e^{-x} \sum_{k=1}^{k^*} \rho^k \leq \frac{\rho}{1 - \rho} e^{-x}.$$

Let $\mathcal{Y}(\boldsymbol{v})$, $\boldsymbol{v} \in \Upsilon$, be a smooth centered random vector process with values in $I\!\!R^q$, where $\Upsilon \subseteq I\!\!R^p$. Let also $\mathcal{Y}(\boldsymbol{v}_0) = 0$ for a fixed point $\boldsymbol{v}_0 \in \Upsilon$. (w.l.g. $\boldsymbol{v}_0 = 0$).

Suppose that $\mathcal{Y}(\boldsymbol{v})$ satisfies for each $\boldsymbol{\gamma} \in I\!\!R^p$ and $\boldsymbol{\alpha} \in I\!\!R^q$ with $\|\boldsymbol{\gamma}\| = \|\boldsymbol{\alpha}\| = 1$

$$\sup_{\boldsymbol{v} \in \Upsilon} \log I\!\!E \exp\Big\{\lambda \boldsymbol{\gamma}^\top \nabla \mathcal{Y}(\boldsymbol{v}) \boldsymbol{\alpha}\Big\} \leq \frac{\nu_0^2 \lambda^2}{2}, \qquad \lambda^2 \leq 2\mathsf{g}^2. \tag{7}$$

We aim to bound the maximum of the norm $\|\mathcal{Y}(\boldsymbol{v})\|$ over a ball

$$\Upsilon_\circ(\mathbf{r}) = \big\{\boldsymbol{v} \in \Upsilon \colon \|\boldsymbol{v} - \boldsymbol{v}_0\| \leq \mathbf{r}\big\}.$$

Condition (7) implies for any $\boldsymbol{v} \in \Upsilon_\circ(\mathbf{r})$ with $\|\boldsymbol{v}\| \leq \mathbf{r}$ and $\|\boldsymbol{\gamma}\| = 1$ in view of $\mathcal{Y}(\boldsymbol{v}_0) = 0$

$$\log I\!\!E \exp\Big\{\frac{\lambda}{\mathbf{r}} \boldsymbol{\gamma}^\top \mathcal{Y}(\boldsymbol{v})\Big\} \leq \frac{\nu_0^2 \lambda^2 \|\boldsymbol{v}\|^2}{2\mathbf{r}^2}, \qquad \lambda^2 \leq 2\mathsf{g}^2; \tag{8}$$

Use the representation

$$\|\mathcal{Y}(\boldsymbol{v})\| = \sup_{\|\boldsymbol{u}\| \leq \mathtt{r}} \frac{1}{\mathtt{r}} \boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v}).$$

This implies for $\varUpsilon_\circ(\mathtt{r}) = \{\boldsymbol{v} \in \varUpsilon \colon \|\boldsymbol{v} - \boldsymbol{v}_0\| \leq \mathtt{r}\}$

$$\sup_{\boldsymbol{v} \in \varUpsilon_\circ(\mathtt{r})} \|\mathcal{Y}(\boldsymbol{v})\| = \sup_{\boldsymbol{v} \in \varUpsilon_\circ(\mathtt{r})} \sup_{\|\boldsymbol{u}\| \leq \mathtt{r}} \frac{1}{\mathtt{r}} \boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v}).$$

Consider a bivariate process $\boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v})$ of $\boldsymbol{u} \in \mathbb{R}^q$ and $\boldsymbol{v} \in \Upsilon \subset \mathbb{R}^p$.
By definition $\mathbb{E}\boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v}) = 0$. Further, for $\boldsymbol{\gamma} = \boldsymbol{u}/\|\boldsymbol{u}\|$

$$\nabla_{\boldsymbol{u}}\big[\boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v})\big] = \mathcal{Y}(\boldsymbol{v}), \qquad \nabla_{\boldsymbol{v}}\big[\boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v})\big] = \boldsymbol{u}^\top \nabla \mathcal{Y}(\boldsymbol{v}) = \|\boldsymbol{u}\|\boldsymbol{\gamma}^\top \nabla \mathcal{Y}(\boldsymbol{v})$$

Suppose that $\boldsymbol{u} \in \mathbb{R}^q$ and $\boldsymbol{v} \in \Upsilon$ are such that $\|\boldsymbol{u}\|^2 + \|\boldsymbol{v}\|^2 \le 2\mathbf{r}^2$. By the Hölder inequality, (8), and (7), it holds for $\|\boldsymbol{\gamma}\| = \|\boldsymbol{\alpha}\| = 1$ and $\boldsymbol{v} \in \Upsilon_\circ(\mathbf{r})$

$$\log \mathbb{E}\exp\bigg\{\frac{\lambda}{2\mathbf{r}}(\boldsymbol{\gamma}, \boldsymbol{\alpha})^\top \nabla\big[\boldsymbol{u}^\top \mathcal{Y}(\boldsymbol{v})\big]\bigg\}$$

$$\le \frac{1}{2}\log \mathbb{E}\exp\bigg\{\frac{\lambda}{\mathbf{r}}\boldsymbol{\gamma}^\top \mathcal{Y}(\boldsymbol{v})\bigg\} + \frac{1}{2}\log \mathbb{E}\exp\bigg\{\frac{\lambda}{\mathbf{r}}\boldsymbol{u}^\top \nabla \mathcal{Y}(\boldsymbol{v})\boldsymbol{\alpha}\bigg\}$$

$$\le \frac{1}{2}\log \mathbb{E}\exp\bigg\{\frac{\lambda}{\mathbf{r}}\boldsymbol{\gamma}^\top \mathcal{Y}(\boldsymbol{v})\bigg\} + \frac{1}{2}\log \mathbb{E}\exp\bigg\{\frac{\lambda}{\mathbf{r}}\|\boldsymbol{u}\|\boldsymbol{\gamma}^\top \nabla \mathcal{Y}(\boldsymbol{v})\boldsymbol{\alpha}\bigg\}$$

$$\le \frac{\nu_0^2 \lambda^2}{4\mathbf{r}^2}\big(\|\boldsymbol{v}\|^2 + \|\boldsymbol{u}\|^2\big) \le \frac{\nu_0^2 \lambda^2}{2}, \qquad |\lambda| \le \mathsf{g}.$$

**Theorem**

*Let a random $p$-vector process $\mathcal{Y}(\boldsymbol{v})$ for $\boldsymbol{v} \in \Upsilon \subseteq I\!\!R^p$ fulfill $\mathcal{Y}(\boldsymbol{v}_0) = 0$, $I\!\!E\mathcal{Y}(\boldsymbol{v}) \equiv 0$, and the condition (7) be satisfied. Then for each $\mathtt{r}$ and any $\mathtt{x} \geq 1/2$, it holds*

$$I\!\!P\Big\{ \sup_{\boldsymbol{v} \in \Upsilon_\circ(\mathtt{r})} \big\|\mathcal{Y}(\boldsymbol{v})\big\| > 6\nu_0 \mathtt{r}\, z_{\mathbb{H}}(\mathtt{x}) \Big\} \leq \mathrm{e}^{-\mathtt{x}},$$

*where $z_{\mathbb{H}}(\mathtt{x})$ is given by the following rule: with $\mathbb{H} = 4p$*

$$z_{\mathbb{H}}(\mathtt{x}) = \begin{cases} \sqrt{\mathbb{H} + 2\mathtt{x}} & \text{if } \mathbb{H} + 2\mathtt{x} \leq \mathtt{g}^2, \\ \mathtt{g}^{-1}\mathtt{x} + \frac{1}{2}\big(\mathtt{g}^{-1}\mathbb{H} + \mathtt{g}\big) & \text{if } \mathbb{H} + 2\mathtt{x} > \mathtt{g}^2, \end{cases}$$

# Outline

**I.i.d. model**

Consider $\boldsymbol{Y} = (Y_1, \ldots, Y_n)^\top$.

PA: $Y_i$ i.i.d. from $P \in (P_{\boldsymbol{\theta}})$ with a log-density $\ell(y, \boldsymbol{\theta})$.

Yields

$$L(\boldsymbol{\theta}) \,=\, \sum_{i=1}^n \ell(Y_i, \boldsymbol{\theta}),$$

$$\widetilde{\boldsymbol{\theta}} \,=\, \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}),$$

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

True: $Y_i$'s are i.i.d. from $P \notin (P_{\boldsymbol{\theta}})$,

$$D_n^2 = n I\!\!F_{\boldsymbol{\theta}^*}.$$

(for simplicity $p = 1$)

▶ Smoothness:

- ■ $\nabla^2 I\!E \ell(Y_1, \boldsymbol{\theta})$ Lipschitz continuous in $\boldsymbol{\theta}$;

- ■ $I\!E \exp\big\{\lambda_0 \ell'(Y_1, \boldsymbol{\theta})\big\} \leq \mathtt{C}$

- ■ $I\!E \exp\big\{\lambda_0 \ell''(Y_1, \boldsymbol{\theta})\big\} \leq \mathtt{C}$

▶ Identifiability:

$-\nabla^2 I\!E \ell(\boldsymbol{\theta}) > 0$ and $\Theta$ compact;

Then the conditions are fulfilled with $\mathtt{g}^2 \approx n\lambda_0$ and $\mathtt{b(r)} \geq \mathtt{b}_0 > 0$.

**Checking** $(ED_0)$

Define $\zeta_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \ell(Y_i, \boldsymbol{\theta}) - I\!E\ell(Y_i, \boldsymbol{\theta})$.

Let

$$\mathbf{v}_0^2 = \text{Var}\{\nabla\zeta_i(\boldsymbol{\theta}^*)\}, \qquad V_0^2 = n\mathbf{v}_0^2$$

and

$$\log \exp\{\lambda n^{-1/2}\mathbf{v}_0^{-1}\nabla\zeta_i(\boldsymbol{\theta}^*)\} \leq \frac{\nu_0^2\lambda^2}{2}, \qquad |\lambda| \leq \mathsf{g}_0$$

Then for $|\lambda| \leq \mathsf{g}_0 n^{1/2}$

$$\log I\!E \exp\{\lambda V_0^{-1}\nabla\zeta(\boldsymbol{\theta}^*)\} = \sum_i \log \exp\{\lambda n^{-1/2}\mathbf{v}_0^{-1}\nabla\zeta_i(\boldsymbol{\theta}^*)\} \leq \frac{\nu_0^2\lambda^2}{2}.$$

Let $\boldsymbol{Y} = (Y_1, \ldots, Y_n)^\top \sim I\!\!P$, a sample of independent r.v.s.

Consider PA: $Y_i \sim P_{\Psi_i^\top \boldsymbol{\theta}} \in (P_{\boldsymbol{v}})$, where

– $\Psi_i$, given factors in $I\!\!R^p$,

– $(P_{\boldsymbol{v}})$, an exponential family with canonical parametrization, $\ell(y, \boldsymbol{v}) = y\boldsymbol{v} - d(\boldsymbol{v})$,

– $\boldsymbol{\theta} \in I\!\!R^p$, unknown parameter.

MLE:

$$\widetilde{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^{n} \big\{ Y_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta}) \big\}$$

$$\boldsymbol{\theta}^* = \operatorname*{argmax}_{\boldsymbol{\theta}} I\!\!E L(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^{n} \big\{ f_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta}) \big\}$$

with $f_i = I\!\!E Y_i$.

$$L(\boldsymbol{\theta}) = \sum_{i=1}^{n} \big\{ Y_i \Psi_i^\top \boldsymbol{\theta} - d(\Psi_i^\top \boldsymbol{\theta}) \big\}.$$

Stochastic component is linear in $\boldsymbol{\theta}$

$$\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - I\!\!EL(\boldsymbol{\theta}) = \Big( \sum_{i=1}^{n} \varepsilon_i \Psi_i \Big)^\top \boldsymbol{\theta}$$

$\nabla^2 \zeta(\boldsymbol{\theta}) \equiv 0$ and $(ED_2)$ automatically;
Fisher information only depends on the model:

$$D_0^2 = \sum_i \Psi_i \Psi_i^\top \, d''(\Psi_i^\top \boldsymbol{\theta}^*)$$

The vector $\boldsymbol{\xi}$:

$$\boldsymbol{\xi} = D_0^{-1} \nabla \zeta(\boldsymbol{\theta}^*) = D_0^{-1} \sum_{i=1}^{n} \varepsilon_i \Psi_i$$

Sufficient conditions:

– $d''(\varPsi_i^\top \boldsymbol{\theta})$ uniformly continuous in $\boldsymbol{\theta}$ over $i = 1, \ldots, n$;

– for some fixed matrices $\mathbf{v}_i^2$ and $\lambda_0 > 0$

$$I\!\!E \exp\{\lambda \mathbf{v}_i^{-1} \varepsilon_i\} \leq \mathtt{C}$$

– the matrix $V_0^2 = \sum_i \mathbf{v}_i^2$ fulfills

$$V_0^2 \leq \mathfrak{a}^2 D_0^2$$

The Fisher expansion is simple because the stochastic term is linear in parameter $\boldsymbol{\theta}$. Only smoothness of $d(\boldsymbol{\upsilon})$ and exponential moments of $Y_i$ are required.

Consider a median linear regression

$$Y_i = \Psi_i^\top \boldsymbol{\theta} + \varepsilon_i, \qquad \text{med}(\varepsilon_i) = 0.$$

PA: $Y_i - \Psi_i^\top \boldsymbol{\theta} \sim$ i.i.d. Laplace . Yields

$$L(\boldsymbol{\theta}) = -\sum_i |Y_i - \Psi_i^\top \boldsymbol{\theta}| + R$$

MLE = LAD

$$\widetilde{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname*{argmin}_{\boldsymbol{\theta}} \sum_i |Y_i - \Psi_i^\top \boldsymbol{\theta}|$$

$$\boldsymbol{\theta}^* = \operatorname*{argmax}_{\boldsymbol{\theta}} I\!\!E L(\boldsymbol{\theta}) = \operatorname*{argmin}_{\boldsymbol{\theta}} \sum_i I\!\!E |Y_i - \Psi_i^\top \boldsymbol{\theta}|$$

Sufficient conditions:

– the density $f_i(0)$ of $\varepsilon_i = Y_i - \Psi_i^\top \boldsymbol{\theta}^*$ satisfy

$$f_i(0) \geq \mathtt{C} > 0$$

the sample size $n$ satisfies

$$n \geq \mathtt{C}p$$

Observed $Z_i = (X_i, Y_i)$.

Conditional estimating equations (or moment restrictions)

$$\mathbb{E}\big[g(Z, \boldsymbol{\theta}) \,\big|\, X\big] = 0 \quad \text{a.s.} \quad \Leftrightarrow \quad \boldsymbol{\theta} = \boldsymbol{\theta}_0.$$

Here $g(Z, \boldsymbol{\theta})$ is a known function, of $Z$ and $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p$.

Common models that fit into this framework are

1. (non)linear regression models: $g(Z, \boldsymbol{\theta}) = Y - f(X, \boldsymbol{\theta})$;

2. conditional quantile models: $g(Z, \boldsymbol{\theta}) = \mathbb{1}\big\{Y - f(X, \boldsymbol{\theta}) \leq 0\big\} - \tau$ for a quantile of order $\tau$;

3. linear transformation regression models: $g(Z, \boldsymbol{\theta}) = h(Y, \boldsymbol{\eta}) - X^\top \boldsymbol{\beta}$ and $\boldsymbol{\theta} = (\eta^\top, \boldsymbol{\beta}^\top)^\top$;

4. instrumental variables models;

5. econometric models of optimizing agents, e.g. the consumption model of Hansen and Singleton (1982).

▶ A classical approach: exploit a finite number of unconditional estimating equations:

$$\mathbb{E}\big[A(X)g(Z, \boldsymbol{\theta}_0)\big] = 0 \quad \text{a.s.}$$

where $A(X)$ is a user-selected matrix function.

▶ Generalized Method of Moments (GMM) (Hansen, 1982): minimize a weighted quadratic form in the empirical analog of the moment conditions.

▶ Qin and Lawless (1994) develop an empirical likelihood type estimator.

▶ Smooth Minimum Distance (SMD) (Lavergne and Patilea, 2010):

$$\mathbb{E}\big[g(Z_1, \boldsymbol{\theta})^\top g(Z_2, \boldsymbol{\theta})\, \omega(X_1 - X_2)\big],$$

where $Z_1$ and $Z_2$ are two independent copies of $Z$, and

$$\omega(x) = \omega_h(x) = K(x/h),$$

where $h$ is a bandwidth and $K$ is a kernel.

Let $\boldsymbol{Z}$ be the observed data. Define

$$M(\boldsymbol{\theta}) = M(\boldsymbol{Z}, \boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i,j=1}^{n} g_i(\boldsymbol{\theta}) g_j(\boldsymbol{\theta}) w_{ij} \,,$$

where

– $g_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} g(Z_i, \boldsymbol{\theta})$ ,

– $w_{ij}$ is the collection of localizing weights: $w_{ij} = N^{-1} K\left(\frac{X_i - X_j}{h}\right)$ and

– $N$ is a normalizing factor which ensures that

$$\sum_j w_{ij} = \frac{1}{N} \sum_j K\left(\frac{X_i - X_j}{h}\right) \approx 1.$$

Simple calculus yields the expectation

$$\mathbb{E}M(\boldsymbol{\theta}) = \sum_{i,j} b_i(\boldsymbol{\theta})b_j(\boldsymbol{\theta})w_{ij} + \sum_i \mathbb{E}\varepsilon_i^2(\boldsymbol{\theta})\,w_{ii}\,, \qquad (9)$$

where $b_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \mathbb{E}g_i(\boldsymbol{\theta})$ and $\varepsilon_i(\boldsymbol{\theta}) \stackrel{\text{def}}{=} g_i(\boldsymbol{\theta}) - \mathbb{E}g_i(\boldsymbol{\theta}) = g_i(\boldsymbol{\theta}) - b_i(\boldsymbol{\theta})\,.$
Under PA, $\boldsymbol{\theta}^*$ minimizes the first sum in (9):

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \sum_{i,j} b_i(\boldsymbol{\theta})b_j(\boldsymbol{\theta})w_{ij}\,.$$

If the variance $\operatorname{Var} g_i(\boldsymbol{\theta}) = \mathbb{E}\varepsilon_i^2(\boldsymbol{\theta})$ is available, one can consider

$$M^c(\boldsymbol{\theta}) \stackrel{\text{def}}{=} M(\boldsymbol{\theta}) - \sum_i \mathbb{E}\varepsilon_i^2(\boldsymbol{\theta})\,w_{ii}\,.$$

Alternatively, one often leaves the cross terms $g_i^2(\boldsymbol{\theta})w_{ii}$ out in the definition of $M(\boldsymbol{\theta})$

$$M^-(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{i,j\,:\,i\neq j} g_i(\boldsymbol{\theta})g_j(\boldsymbol{\theta})w_{ij}\,.$$

Consider

$$\widetilde{\boldsymbol{\theta}} \overset{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} M^c(\boldsymbol{\theta}) = \sum_{i,j} g_i(\boldsymbol{\theta}) g_j(\boldsymbol{\theta}) w_{ij} - \sum_i I\!E \varepsilon_i^2(\boldsymbol{\theta}) \, w_{ii} \, .$$

Define also

$$\boldsymbol{\theta}^* \overset{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \, I\!E M^c(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \sum_{i,j} b_i(\boldsymbol{\theta}) b_j(\boldsymbol{\theta}) w_{ij}.$$

Under PA $\boldsymbol{b}(\boldsymbol{\theta}^*) \equiv 0$ and $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$ .

Aim: accuracy (root-n consistency, efficiency) of $\widetilde{\boldsymbol{\theta}}$ .

Problem: the quadratic term $\sum_{i,j} \varepsilon_i(\boldsymbol{\theta}) \varepsilon_j(\boldsymbol{\theta}) w_{ij}$ is not sufficiently regular in $\boldsymbol{\theta}$ .

Represent

$$M^c(\boldsymbol{\theta}) = \boldsymbol{g}(\boldsymbol{\theta})^\top W \boldsymbol{g}(\boldsymbol{\theta}) - S(\boldsymbol{\theta}) = \|A\boldsymbol{g}(\boldsymbol{\theta})\|^2 - S(\boldsymbol{\theta})$$

where $AA^\top = W$ and $S(\boldsymbol{\theta}) = \sum_i I\!E \varepsilon_i^2(\boldsymbol{\theta})\, w_{ii}$.

For simplicity suppose that $S(\boldsymbol{\theta})$ is smooth or constant in the vicinity of $\boldsymbol{\theta}^*$.

Define

$$\boldsymbol{g}_0(\boldsymbol{\theta}) = \boldsymbol{b}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}(\boldsymbol{\theta}^*).$$

Obviously, with $\boldsymbol{g}(\boldsymbol{\theta}) = \boldsymbol{b}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}(\boldsymbol{\theta})$, it holds

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\{ \|A\boldsymbol{g}(\boldsymbol{\theta})\| - \|A\boldsymbol{g}_0(\boldsymbol{\theta})\| \right\} \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \left\| A\{\boldsymbol{\varepsilon}(\boldsymbol{\theta}) - \boldsymbol{\varepsilon}(\boldsymbol{\theta}^*)\} \right\|.$$

Idea: consider separately $\|A\boldsymbol{g}_0(\boldsymbol{\theta})\|^2$ and $\left\| A\{\boldsymbol{\varepsilon}(\boldsymbol{\theta}) - \boldsymbol{\varepsilon}(\boldsymbol{\theta}^*)\} \right\|$.

### Theorem

*Suppose that for any $\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)$ and each $i = 1, \ldots, n$ and any unit vector $\boldsymbol{\gamma} \in \mathbb{R}^p$*

$$\log \mathbb{E} \exp\left\{\lambda \boldsymbol{\gamma}^\top \nabla \varepsilon_i(\boldsymbol{\theta})\right\} \leq \nu_0^2 \lambda^2 / 2, \qquad \lambda^2 \leq 2\mathbf{g}^2,$$

*Then for each $\mathbf{r}$, it holds on a random set $\Omega_1(\mathbf{x})$ of a dominating probability at least $1 - \mathrm{e}^{-\mathbf{x}}$*

$$\sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r})} \|A\varepsilon(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\| \leq 6\nu_0 \, \mathbf{r} \, \mathfrak{z}_A(\mathbf{x}) / \sqrt{N}.$$

*where the function $\mathfrak{z}_A(\mathbf{x})$ is given by*

$$\mathfrak{z}_A(\mathbf{x}) = \mathbb{H}_1 + \sqrt{2\mathbf{x}} + \mathbf{g}^{-1}(\mathbf{g}^{-2}\mathbf{x} + 1)\mathbb{H}_2.$$

*Here $\mathbb{H}_1 = 2\mathbb{H}_1(A)$ and $\mathbb{H}_2 = \mathbb{H}_2(A) + 2\mathfrak{c}_1 p$ with*

$$\mathbb{H}_2(A) = 1 + \frac{8}{3}\operatorname{tr}(A^{-1}), \quad \mathbb{H}_1(A) = 1 + 2\sqrt{\operatorname{tr}(A^{-2}\log(A^2))}.$$

he major step in our study is a local linear approximation of $L_0(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\|A\boldsymbol{g}_0(\boldsymbol{\theta})\|^2/2$:

$$L_0(\boldsymbol{\theta}) \stackrel{\text{def}}{=} -\frac{1}{2}\|A\boldsymbol{g}_0(\boldsymbol{\theta})\|^2 = -\frac{1}{2}\big\|A\big\{\boldsymbol{b}(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}(\boldsymbol{\theta}^*)\big\}\big\|^2.$$

It is obvious that

$$I\!\!E L_0(\boldsymbol{\theta}) = -\frac{1}{2}\|A\boldsymbol{b}(\boldsymbol{\theta})\|^2 - \frac{1}{2}I\!\!E\|A\boldsymbol{\varepsilon}(\boldsymbol{\theta}^*)\|^2.$$

This implies that $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} I\!\!E L_0(\boldsymbol{\theta})$. Further, define

$$D_0^2 = -\nabla^2 I\!\!E L_0(\boldsymbol{\theta}^*) = -\frac{1}{2}\sum_{i,j}\big\{\nabla^2 b_i b_j\big\}(\boldsymbol{\theta}^*)w_{ij} = -\frac{1}{2}\nabla^2\big\{\boldsymbol{b}^\top W \boldsymbol{b}\big\}(\boldsymbol{\theta}^*).$$

In the our case when PA is correct, it holds $b_i(\boldsymbol{\theta}^*) \equiv 0$ and

$$D_0^2 = -\frac{1}{2}\sum_{i,j}\nabla b_i(\boldsymbol{\theta}^*)\big\{\nabla b_j(\boldsymbol{\theta}^*)\big\}^\top w_{ij} = -\nabla\boldsymbol{b}(\boldsymbol{\theta}^*)^\top W \nabla\boldsymbol{b}(\boldsymbol{\theta}^*).$$

# Outline

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

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

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

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

Introduce the posterior moments

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

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

## 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;

**Theorem**

*Suppose the conditions of Theorem 25. Let also* $\mathtt{b}(\mathbf{r})$ *from* $(\mathcal{L})$ *satisfies*

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

*with* $\varrho(\mathbf{r}, \mathbf{x})$ *from* (22). *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 (10) 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

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

**Theorem**

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

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

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

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

**Theorem**

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

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

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

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

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

► All statements of Theorem 13 require " $\Delta(\mathbf{r}_0, \mathbf{x})$ 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 $\|D_0(\widehat{\boldsymbol{\theta}} - \breve{\boldsymbol{\theta}})\| \leq \mathbf{r}_0$ with a dominating probability.

**Credible sets**

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

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

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 $1 - 5\mathrm{e}^{-\mathtt{x}}$*

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

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

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

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

Now the result follows from Theorem 13 and

**Lemma**

*Let* $\mathbb{P}_0 = \mathcal{N}(0, \boldsymbol{I}_p)$ *and* $\mathbb{P}_1 = \mathcal{N}(\boldsymbol{\beta}, (U^\top U)^{-1})$ *some non-degenerated matrix* $U$. *If*

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

*then* $\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) = -\mathbb{E}_0 \log \frac{d\mathbb{P}_1}{d\mathbb{P}_0}$ *fulfills*

$$2\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1) \;\leq\; \operatorname{tr}(U^\top U - \boldsymbol{I}_p)^2 + (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2 \leq \boldsymbol{\epsilon}^2\, p + (1 + \boldsymbol{\epsilon})\|\boldsymbol{\beta}\|^2.$$

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

$$\big|\mathbb{P}_0(A) - \mathbb{P}_1(A)\big| = \big|\mathbb{P}(\boldsymbol{\gamma} \in A) - \mathbb{P}(U(\boldsymbol{\gamma} - \boldsymbol{\beta}) \in A)\big| \leq \sqrt{\mathcal{K}(\mathbb{P}_0, \mathbb{P}_1)/2}.$$

**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

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

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

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

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

This implies by Pinsker's inequality

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

Remind

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

and $\boldsymbol{\xi} = D_0^{-1}\nabla L(\boldsymbol{\theta}^*)$. For any nonnegative function $f$, it holds

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

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

Similarly,

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

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

The main benefit of these bounds is that $\mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$ is quadratic in $\boldsymbol{\theta}$.

## Theorem

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

$$\mathbb{E}^\circ\left[f\left(D_0(\boldsymbol{\vartheta} - \check{\boldsymbol{\theta}})\right)\mathbb{1}_{\mathbf{r}_0}\right] \leq \exp\{\Delta^+(\mathbf{r}_0, \mathbf{x})\}\,\mathbb{E}f(\boldsymbol{\gamma}), \tag{11}$$

*where*

$$\Delta^+(\mathbf{r}_0, \mathbf{x}) = 2\Delta(\mathbf{r}_0, \mathbf{x}) + \nu(\mathbf{r}_0), \tag{12}$$

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

*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

$$m(\boldsymbol{\xi}) + \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) \;=\; -\|D_0(\boldsymbol{\theta} - \breve{\boldsymbol{\theta}})\|^2/2 + \log(\det D_0) - p\log(\sqrt{2\pi}) \tag{13}$$

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 (13) for any nonnegative function $f$ that

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

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

$$= \; \mathrm{e}^{\Delta(\mathbf{r}_0, \mathbf{x})} \int \phi(\boldsymbol{u}) \, f(\boldsymbol{u}) \, d\boldsymbol{u} = \mathrm{e}^{\Delta(\mathbf{r}_0, \mathbf{x})} \, I\!\!E f(\boldsymbol{\gamma}). \tag{14}$$

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

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

$$\geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi})\} \int \phi(\boldsymbol{u}) f(\boldsymbol{u}) \, \mathbb{1}\big\{\|\boldsymbol{u} + \boldsymbol{\xi}\| \leq \mathbf{r}_0\big\} d\boldsymbol{u}. \tag{15}$$

A special case of (15) with $f(\boldsymbol{u}) \equiv 1$ implies by definition of $\nu(\mathbf{r}_0)$ :

$$\int_{\Theta_0(\mathbf{r}_0)} \exp\{L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)\} \, d\boldsymbol{\theta} \geq \exp\{-\Delta(\mathbf{r}_0, \mathbf{x}) - m(\boldsymbol{\xi}) - \nu(\mathbf{r}_0)\}. \tag{16}$$

Now we are prepared to finalize the proof. (14) and (16) 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 (11) follows. As $\|\boldsymbol{\xi}\| \leq z(B, \mathbf{x})$ on $\Omega(B, \mathbf{x})$ and $\mathbf{r}_0 \geq z(B, \mathbf{x}) + z(p, \mathbf{x})$,

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

**Lemma**

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

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

*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

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

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 $1 - \mathrm{e}^{-\mathbf{x}}$ such that

$$\big|\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\big| \leq \mathbf{r}\,\varrho(\mathbf{r}, \mathbf{x}) \tag{17}$$

for $\mathbf{r} = \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|$ and $\varrho(\mathbf{r}, \mathbf{x})$ from (4); cf. the proof of Theorem 25.
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*

$$\mathbf{r}^2\mathbf{b}^2(\mathbf{r}) \;\geq\; \mathbf{x} + 2p + 4z^2(B,\mathbf{x}) + 8\mathbf{r}\,\mathbf{b}(\mathbf{r})\varrho(\mathbf{r},\mathbf{x}), \qquad \mathbf{r} \geq \mathbf{r}_0, \qquad (18)$$

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

$$\rho(\mathbf{r}_0) \;\overset{\text{def}}{=}\; \frac{\int_{\Theta\setminus\Theta_0} \exp\{L(\boldsymbol{\theta})\}d\boldsymbol{\theta}}{\int_{\Theta_0} \exp\{L(\boldsymbol{\theta})\}d\boldsymbol{\theta}} \;\leq\; 2\exp\{-\mathbf{x} + \Delta^+(\mathbf{r}_0,\mathbf{x})\} \qquad (19)$$

*with* $\Delta^+(\mathbf{r}_0,\mathbf{x})$ *from* (12). *Further, for any unit vector* $\boldsymbol{a} \in \mathbb{R}^p$

$$\rho_2(\mathbf{r}_0) \;\overset{\text{def}}{=}\; \frac{\int_{\Theta\setminus\Theta_0} |\boldsymbol{a}^\top D_0(\boldsymbol{\theta}-\boldsymbol{\theta}^*)|^2 \exp\{L(\boldsymbol{\theta},\boldsymbol{\theta}^*)\}d\boldsymbol{\theta}}{\int_{\Theta_0} \exp\{L(\boldsymbol{\theta},\boldsymbol{\theta}^*)\}d\boldsymbol{\theta}} \;\leq\; 2\exp\{-\mathbf{x} + \Delta^+(\mathbf{r}_0,\mathbf{x})\}.$$

Suppose that $\mathbf{b}_0 = \mathbf{b}(\mathbf{r}_0)$ is close to one. Condition (18) 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

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

for each $k \geq 0$ and any $\boldsymbol{\theta} \in \Theta_0(\mathtt{r}_{k+1}) \setminus \Theta_0(\mathtt{r}_k)$. The bound (17) implies on $\Omega(\mathbf{x})$

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

for all $k \geq 0$.

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

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

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

Next,

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

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

$$\leq \mathtt{b}_k^{-p/2} \exp\left(\frac{\|\boldsymbol{\xi}\|^2}{\mathtt{b}_k} - \frac{1}{4}\mathtt{b}_k\mathtt{r}_k^2 + \frac{p}{2}\right). \tag{20}$$

Here we have used the bound for a standard normal vector $\boldsymbol{\gamma}$ and $\boldsymbol{u} = \mathtt{b}_k^{-1/2}\boldsymbol{\xi} \in I\!\!R^p$. (16) and (20) imply (19).

Now the bound $\|\boldsymbol{\xi}\| \le z(B, \mathtt{x})$ holding with a dominating probability and (18) imply

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

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

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

and (19) follows in view of $\mathtt{b}\log(e/\mathtt{b}) \le 1$ for $\mathtt{b} \le 1$.

**Theorem**

*Suppose* (**??**) *for* $\mathbf{r} = \mathbf{r}_0$ *and* (19). *Then for any nonnegative function* $f(\cdot)$ *on* $\mathbb{R}^p$, *it holds on* $\Omega(\mathbf{x})$

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

*where*

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

## Lower bound. Proof

On the set $\Omega(\mathbf{x})$, it holds by (14) 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 (15) imply

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

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

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

# Outline

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

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

■ $p_n/n \to 0$ : Consistency:

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

■ $p_n^2/n \to 0$ – 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.

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.

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 $\|G\boldsymbol{\theta}_G^*\|^2$.

"Bias-variance" trade-off:

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

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

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

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

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

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

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

$$\mathtt{p}_G \stackrel{\mathrm{def}}{=} \mathrm{tr}\big(B_G\big).$$

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

$$D_G^2 = D_0^2 + G^2 = \text{diag}\{\sigma^2 + g_1^2, \ldots, \sigma^2 + g_p^2\},$$

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

## Block penalization

$G$ 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

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

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

For $\beta > 1/2$,

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

$$g_j = L j^\beta$$

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

It holds

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

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

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

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

**Theorem**

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

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

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

*with*

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

$$\widetilde{\boldsymbol{\theta}}_G \overset{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta}{\text{argmax}}\, L_G(\boldsymbol{\theta}), \qquad \boldsymbol{\theta}_G^* \overset{\text{def}}{=} \underset{\boldsymbol{\theta} \in \Theta}{\text{argmax}}\, I\!\!E L_G(\boldsymbol{\theta})$$

**Theorem**

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

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

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

*with*

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

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

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

$$\varrho(\mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} 6\nu_0\, z_{\mathbb{H}}\big(\mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0)\big)\,\omega. \tag{21}$$

*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*

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

---

**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*

$$\varrho(\mathbf{r}, \mathbf{x}) \overset{\text{def}}{=} 6\nu_0 \, z_{\mathbb{H}}\big( \mathbf{x} + \log(2\mathbf{r}/\mathbf{r}_0) \big) \, \omega. \tag{22}$$

*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 $\Upsilon_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{v} \colon \|\boldsymbol{v}\| \le \mathbf{r}\}$

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

---

**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 $\Upsilon_\circ(\mathbf{r}) \stackrel{\text{def}}{=} \{\boldsymbol{v} \colon \|B^{-1/2}\boldsymbol{v}\| \leq \mathbf{r}\}$

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

---

**Theorem**

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

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

*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*

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

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

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

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

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

---

**Theorem**

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

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

*where*

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

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 $\boldsymbol{\theta} \in \Theta_{0,G}(\mathbf{r})$

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

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

---

### 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

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

- $p_n/n \to 0$ : Consistency:

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

The semiparametric Bernstein-von Mises theorem.

*Ann. Statist.*, 40(1):206–237.

Birgé, L. and Massart, P. (1993).

Rates of convergence for minimum contrast estimators.

*Probab. Theory Relat. Fields*, 97(1-2):113–150.

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.

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.

📄 Portnoy, S. (1984).

Asymptotic behavior of M-estimators of p regression parameters when $p^2/n$ is large. I. Consistency.

*Ann. Stat.*, 12:1298–1309.

📄 Portnoy, S. (1988).

Asymptotic behavior of likelihood methods for exponential families when the number of parameters tends to infinity.

*Ann. Statist.*, 16(1):356–366.

📄 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 de Geer, S. (1993).

Hellinger-consistency of certain nonparametric maximum likelihood estimators.

*Ann. Stat.*, 21(1):14–44.

📄 van de Geer, S. (2002).

M-estimation using penalties or sieves.

*J. Stat. Plann. Inference*, 108(1-2):55–69.