



Weierstraß-Institut für  
Angewandte Analysis und Stochastik



# Fisher and Wilks expansions with applications to statistical inference

Vladimir Spokoiny,  
WIAS, HU Berlin

### 1 Introduction. Fisher and Wilks expansions

- Fisher and Wilks expansions
- The case of a linear model
- Expansions vs asymptotic results

### 2 Fisher and Wilks: Main steps

- Local quadraticity of  $\mathcal{EL}(\theta)$
- Local linear approximation of the stochastic term
- Local linear approximation of the gradient and the “Fisher” trick
- Local quadratic approximation of the log-likelihood and the “Wilks” trick
- Concentration and large deviation for  $\tilde{\theta}$
- A sharp bound for  $\|\xi\|^2$
- An upper function for the stochastic component

### 3 Examples

- Summary
- An i.i.d. case
- Generalized linear models
- Linear median (quantile) regression
- Conditional Moment Restriction (CMR)

Data  $\mathbf{Y} \sim \mathbb{P}$ . Aim: infer on  $\mathbb{P}$ .

Parametric assumption (PA):  $\mathbb{P} \in (\mathbb{P}_\theta, \theta \in \Theta \subseteq \mathbb{R}^p) \ll \mu_0$ .

Maximum likelihood estimator (MLE):

$$\tilde{\theta} \stackrel{\text{def}}{=} \operatorname{argmax}_{\theta \in \Theta} L(\theta) = \operatorname{argmax}_{\theta \in \Theta} \log \frac{d\mathbb{P}_\theta}{d\mu_0}(\mathbf{Y})$$

PA-PW:  $\mathbb{P} \notin (\mathbb{P}_\theta)$ . Target of estimation ?

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

Under PA:  $\mathbb{P} = \mathbb{P}_{\theta^*}$  and

$$\operatorname{argmin}_{\theta \in \Theta} \mathcal{K}(\mathbb{P}_{\theta^*}, \mathbb{P}_\theta) = \theta^*$$

$$\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}), \quad \boldsymbol{\theta}^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E} L(\boldsymbol{\theta})$$

### Theorem

On a set  $\Omega(x)$  with  $\mathbb{P}(\Omega(x)) \geq 1 - Ce^{-x}$

$$\|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq \diamond(x),$$

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

with

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

Here  $\diamond(x)$  and  $\Delta(x)$  are explicit error terms.

Given

- $\mathbf{Y}$ , response,
- $\Sigma = \text{Cov}(\mathbf{Y})$ , its covariance matrix
- $\Psi$ , design matrix of regressors:

$$\mathbf{Y} = \Psi^\top \boldsymbol{\theta} + \boldsymbol{\varepsilon}, \quad \mathbb{E}\boldsymbol{\varepsilon} = 0, \quad \text{Cov}(\boldsymbol{\varepsilon}) = \Sigma.$$

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

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

Study under true:  $\mathbb{E}\mathbf{Y} = \mathbf{f}$  and  $\text{Cov}(\mathbf{Y}) = \Sigma_0$ .

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

### Lemma

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

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

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

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

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

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

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

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

## Fisher and Wilks for a linear model

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

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

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

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

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

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

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

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

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

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

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

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

$$\text{Var}(\boldsymbol{\xi}) = \text{Var}(D_0^{-1} \boldsymbol{\Psi} \Sigma^{-1} \boldsymbol{\varepsilon}) = D_0^{-1} \boldsymbol{\Psi} \Sigma^{-1} \text{Var}(\boldsymbol{\varepsilon}) \Sigma^{-1} \boldsymbol{\Psi} D_0^{-1} = I_p.$$

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

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

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

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

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

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

- ▶ 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  $\xi$** .
- ▶ **PA-PW:** Let  $\text{Var}(\mathbf{Y}) = \Sigma_0 \neq \Sigma$ . Then with  $D_0^2 = \Psi \Sigma^{-1} \Psi^\top$

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

and (the **sandwich formula**)

$$\text{Var}(\boldsymbol{\xi}) = \text{Var}\{D_0^{-1} \nabla L(\boldsymbol{\theta}^*)\} = D_0^{-1} V_0^2 D_0^{-1} \neq I_p.$$

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

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

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \ell(Y_i, \boldsymbol{\theta}), \quad \tilde{\boldsymbol{\theta}}_n = \operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}).$$

### Theorem

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

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

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

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:

$$\|D_0(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_n\| \leq \diamond(\mathbf{x}),$$

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

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  $\text{Var}\{\nabla \ell(Y_1, \boldsymbol{\theta}^*)\} = \mathbb{F}_{\boldsymbol{\theta}^*}$ , and by CLT

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

For

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

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

$$\|D_n(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}^*) - \boldsymbol{\xi}_n\| \leq \diamond_n(\mathbf{x}),$$

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

The error terms satisfy

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

and

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

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

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

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

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

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

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

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

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

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

[?]: M-estimator i.i.d. or linear models:

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

[?]: MLE for a GLM:

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

Sieve estimation:

[?], [Chen, 1993, 1997], [?, ?]; ...

### 1 Introduction. Fisher and Wilks expansions

- Fisher and Wilks expansions
- The case of a linear model
- Expansions vs asymptotic results

### 2 Fisher and Wilks: Main steps

- Local quadraticity of  $\mathbb{E}L(\theta)$
- Local linear approximation of the stochastic term
- Local linear approximation of the gradient and the “Fisher” trick
- Local quadratic approximation of the log-likelihood and the “Wilks” trick
- Concentration and large deviation for  $\tilde{\theta}$
- A sharp bound for  $\|\xi\|^2$
- An upper function for the stochastic component

### 3 Examples

- Summary
- An i.i.d. case
- Generalized linear models
- Linear median (quantile) regression
- Conditional Moment Restriction (CMR)

Aim:

- minimal non-restrictive and natural conditions
- possibly sharp bounds
- all constants explicit, no asymptotic arguments
- model misspecification incorporated
- self-contained

## Main steps

---

- Concentration and large deviations: for some  $r_0$

$$\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(r_0)) \leq e^{-x},$$

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

- Local quadratic approximation of the expected log-likelihood:

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

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

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

- Overall error of the Fisher expansion  $r_0\{\delta(r_0) + \varrho(r_0, x)\}$ ,  
of the Wilks  $r_0^2\{\delta(r_0) + \varrho(r_0, x)\}$ .

Define

$$D^2(\theta) \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\theta).$$

Then  $D_0^2 = D^2(\theta^*)$ .

( $\mathcal{L}_0$ ) For each  $r \leq r_0$ , there is a constant  $\delta(r) \leq 1/2$  such that it holds for any  $\theta \in \Theta_0(r) = \{\theta \in \Theta : \|D_0(\theta - \theta^*)\| \leq r\}$ :

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

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

$$|-2\mathbb{E}L(\theta) + 2\mathbb{E}L(\theta^*) - \|D_0(\theta - \theta^*)\|^2| \leq \delta(r)r^2,$$

$$\begin{aligned} & \|D_0^{-1} \{ \nabla \mathbb{E}L(\theta) - \nabla \mathbb{E}L(\theta^*) \} + D_0(\theta - \theta^*) \| \\ & \leq \| \{ I_p - D_0^{-1} D^2(\theta^*) D_0^{-1} \} D_0(\theta - \theta^*) \| \leq \delta(r)r. \end{aligned}$$

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

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

**(ED<sub>2</sub>)** There exist a value  $\omega > 0$  and for each  $\mathbf{r} > 0$ , a constant  $g(\mathbf{r}) > 0$  such that  $\zeta(\boldsymbol{\theta}) \stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta})$  satisfies for any  $\boldsymbol{\theta} \in \Theta_0(\mathbf{r})$  :

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

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

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

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

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

$$\sup_{\boldsymbol{\theta} \in \Theta_0(r)} \|D_0^{-1} \{ \nabla \zeta(\boldsymbol{\theta}) - \nabla \zeta(\boldsymbol{\theta}^*) \}\| = \omega \sup_{\mathbf{v} \in \Upsilon_0(r)} \|\mathcal{Y}(\mathbf{v})\|,$$
$$\Upsilon_0(r) \stackrel{\text{def}}{=} \{\mathbf{v}: \|\mathbf{v}\| \leq r\}.$$

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

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

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

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

$$\sup_{\gamma_1, \gamma_2 \in \mathbb{R}^p : \|\gamma_1\| = \|\gamma_2\| = 1} \log \mathbb{E} \exp \left\{ \lambda \gamma_1^\top \nabla \mathcal{Y}(\mathbf{v}) \gamma_2 \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g(\mathbf{r}).$$

### Theorem

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

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

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

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

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

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

$$\|D_0^{-1}\{\nabla IEL(\boldsymbol{\theta}) - \nabla IEL(\boldsymbol{\theta}^*)\} + D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \delta(\mathbf{r})\mathbf{r},$$

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

### Theorem

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

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

where

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

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

## “Fisher” trick

---

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

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

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

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

The use of  $\nabla L(\tilde{\boldsymbol{\theta}}) = 0$  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

$$\begin{aligned}\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) &\stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) + \frac{1}{2} \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2 \\ &= L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ), \quad \boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{x})\end{aligned}$$

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

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

$$\begin{aligned}\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) &\stackrel{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^\circ) - (\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)^\top \nabla L(\boldsymbol{\theta}^\circ) + \frac{1}{2} \|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^\circ)\|^2 \\ &= L(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) - \mathbb{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ), \quad \boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})\end{aligned}$$

### 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  $1 - e^{-x}$ , it holds with any  $\boldsymbol{\theta}, \boldsymbol{\theta}^\circ \in \Theta_0(\mathbf{r})$

$$\frac{|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} \leq \diamond(\mathbf{r}, \mathbf{x}), \quad |\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)| \leq \mathbf{r} \diamond(\mathbf{r}, \mathbf{x}),$$

$$\frac{|\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta})|}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} \leq \diamond(\mathbf{r}, \mathbf{x}), \quad |\alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta})| \leq \mathbf{r} \diamond(\mathbf{r}, \mathbf{x}).$$

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

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

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

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

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

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

Together with (1), this yields

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

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

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

$$\begin{aligned} \left| \left\{ 2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \right\}^{1/2} - \|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\| \right| &\leq \frac{|2L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) - \|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|^2|}{\|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|} \\ &\leq \frac{2|\alpha(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)|}{\|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|} \leq \sup_{\boldsymbol{\theta} \in \Theta_0(\mathbf{r}_0)} \frac{2|\alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^*)|}{\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\|} \leq 2\Diamond(\mathbf{r}_0, \mathbf{x}). \end{aligned}$$

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

## Large deviations for $\tilde{\theta}$ . Main steps

Aim: find  $r_0$  ensuring

$$\mathbb{P}(\tilde{\theta} \notin \Theta_0(r_0)) \leq ce^{-x}.$$

- By definition  $\sup_{\theta \in \Theta_0(r_0)} L(\theta, \theta^*) \geq 0$ . Suffices to check that

$$L(\theta, \theta^*) < 0 \quad \forall \theta \in \Theta \setminus \Theta_0(r_0)$$

- Use the decomposition

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

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

$$\sup_{\theta \in \Theta \setminus \Theta_0(r_0)} \{ \zeta(\theta, \theta^*) - (\theta - \theta^*)^\top \nabla \zeta(\theta^*) - f(\theta, \theta^*) \} \leq 0 \quad \text{w.h.p.}$$

(L) For each  $r$ , there exists  $b(r) > 0$  such that  $rb(r) \rightarrow \infty$  as  $r \rightarrow \infty$  and

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

### Theorem

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

$$b(r)r \geq 2z(B, x) + 2\rho(r, x), \quad r > r_0,$$

where

$$\rho(r, x) \stackrel{\text{def}}{=} 6\nu_0 z_{\mathbb{H}}(x + \log(2r/r_0)) \omega.$$

Then

$$\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(r_0)) \leq 3e^{-x}.$$

The radius  $r_0$  has to fulfill

$$b(r)r \geq 2z(B, x) + 2\varrho(r, x), \quad r > r_0,$$

One can use that

- ▶  $b(r_0) \geq 1 - \delta(r_0) \approx 1$ ,
- ▶ the constant  $\omega$  and thus,  $\varrho(r, x)$ , is small, and
- ▶  $rb(r)$  grows with  $r$ .

A simple rule  $r_0 \geq (2 + \delta)z(B, x)$  for some  $\delta > 0$  works in most of cases.

## A sharp bound for a norm of a stochastic vector $\xi = D_0^{-1} \nabla L(\theta^*)$

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

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

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

$$\|\xi\|^2 = \eta^\top B \eta$$

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

$$p_B \stackrel{\text{def}}{=} \text{tr}(B), \quad v_B^2 \stackrel{\text{def}}{=} 2 \text{tr}(B^2), \quad \lambda_B \stackrel{\text{def}}{=} \lambda_{\max}(B).$$

Note that  $p_B = I\!\!E \|\xi\|^2$ . Moreover, if  $\xi$  is a Gaussian vector then  $v_B^2 = \text{Var}(\|\xi\|^2)$ . If  $V_0^2 = D_0^2$ , then  $\lambda_B = 1$ .

## A bound for the norm $\|\xi\|$ . Cont.

Define  $\mu_c = 2/3$ ,  $p_B = \text{tr}(B)$ ,  $v_B^2 = 2 \text{tr}(B^2)$ , and  $\lambda_B = \lambda_{\max}(B)$

$$g_c \stackrel{\text{def}}{=} \sqrt{g^2 - \mu_c p_B},$$

$$2x_c \stackrel{\text{def}}{=} (g^2/\mu_c - p_B)/\lambda_B + \log \det(I_p - \mu_c B/\lambda_B). \quad (2)$$

### Theorem (SP2012)

Let  $(ED_0)$  hold with  $\nu_0 = 1$  and  $g^2 \geq 2p_B$ . Then for each  $x > 0$

$$P(\|\xi\| \geq z(B, x)) = P(\|B^{1/2}\eta\| \geq z(B, x)) \leq 2e^{-x} + 8.4e^{-x_c},$$

where  $z(B, x)$  is defined with  $y_c^2 \leq p_B + 6\lambda_B x_c$  by

$$z^2(B, x) \stackrel{\text{def}}{=} \begin{cases} p_B + 2v_B x^{1/2}, & x \leq v_B/(18\lambda_B), \\ p_B + 6\lambda_B x, & v_B/(18\lambda_B) < x \leq x_c, \\ |y_c + 2\lambda_B(x - x_c)/g_c|^2, & x > x_c. \end{cases}$$

$$p_B = \text{tr}(B), \quad v_B^2 = 2 \text{tr}(B^2), \quad \lambda_B = \lambda_{\max}(B).$$

$$z^2(B, x) \stackrel{\text{def}}{=} \begin{cases} p_B + 2v_B x^{1/2}, & x \leq v_B/(18\lambda_B), \\ p_B + 6\lambda_B x, & v_B/(18\lambda_B) < x \leq x_c, \\ |y_c + 2\lambda_B(x - x_c)/g_c|^2, & x > x_c. \end{cases}$$

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

- The sub-Gaussian regime for  $x \leq v_B/(18\lambda_B)$
- The Poissonian regime for  $x \leq x_c$
- The value  $x_c$  from (2) is of order  $g^2$ . In all our results we suppose that  $g^2$  and hence,  $x_c$  is sufficiently large;

The quadratic form  $\|\xi\|^2$  can be bounded with a dominating probability by  $p_B + 6\lambda_B x$  for a proper  $x$ .

## A “squared norm” trick

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

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

For simplicity we take here  $B = 1$ .

Aim: to bound  $\|\xi\|^2$ .

A sup-representation:

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

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

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

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

## An upper function for the stochastic component $\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta})$

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

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

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  $1 - e^{-\mathbf{x}}$  that for all  $\mathbf{r} \geq \mathbf{r}_0$  and all  $\boldsymbol{\theta}$  with  $\|D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)\| \leq \mathbf{r}$

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

where

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

The use of  $\nabla \mathbb{E}L(\boldsymbol{\theta}^*) = 0$  yields

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

## Concentration and large deviations. Proof

By definition  $\sup_{\theta \in \Theta_0(r_0)} L(\theta, \theta^*) \geq 0$ . So, it suffices to check that  $L(\theta, \theta^*) < 0$  for all  $\theta \in \Theta \setminus \Theta_0(r_0)$ .

We know

$$\sup_{\theta \in \Theta_0(r)} |L(\theta, \theta^*) - IEL(\theta, \theta^*) - (\theta - \theta^*)^\top \nabla L(\theta^*)| \leq \varrho(r, x) r.$$

Also  $\|\xi\| \leq z(B, x)$  on  $\Omega(x)$  and for each  $r \geq r_0$

$$\begin{aligned} & \sup_{\theta \in \Theta_0(r)} |(\theta - \theta^*)^\top \nabla L(\theta^*)| \\ & \leq \sup_{\theta \in \Theta_0(r)} \|D_0(\theta - \theta^*)\| \times \|D_0^{-1} \nabla L(\theta^*)\| = r \|\xi\| \leq z(B, x) r. \end{aligned}$$

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

$$rb(r) \geq 2z(B, x) + 2\varrho(r, x), \quad r > r_0,$$

ensure  $L(\theta, \theta^*) < 0$  for all  $\theta \notin \Theta_0(r_0)$  with a dominating probability.

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

(ED) There exist  $g > 0$ ,  $\nu_0 \geq 1$ , and a symmetric  $H_0 \geq 0$  s.t. it holds

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

We consider the local sets of the elliptic form  $\mathcal{Y}_o(\mathbf{r}) \stackrel{\text{def}}{=} \{\mathbf{v} : \|H_0(\mathbf{v} - \mathbf{v}^*)\| \leq \mathbf{r}\}$ .

### Theorem

Let (ED) hold with some  $g > 0$ , and a matrix  $H_0$ . For any  $x \geq 0$  and any  $r > 0$

$$I\!\!P \left\{ \sup_{\mathbf{v} \in \mathcal{Y}_o(\mathbf{r})} |\mathcal{U}(\mathbf{v}) - \mathcal{U}(\mathbf{v}^*)| \geq 3\nu_0 \mathbf{r} z_{\mathbb{H}}(x) \right\} \leq e^{-x},$$

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

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

## Tools. An “upper function” device

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

$$|\mathcal{U}(v, v^*)| \leq 3\nu_0 r z_{\mathbb{H}}(x).$$

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

### Theorem

Let  $(ED)$  hold on  $\mathcal{B}_{r^*}(v^*)$ . Given  $r_0 < r^*$ , define  $f(r, r_0)$  for some  $\rho < 1$  as

$$f(r, r_0) = 3\nu_0 r z_{\mathbb{H}}(x + \log(r/r_0)), \quad r_0 \leq r \leq r^*. \quad (3)$$

Then it holds

$$\mathbb{P} \left( \sup_{r_0 \leq r \leq r^*} \sup_{v \in \mathcal{Y}_o(r)} \{ \mathcal{U}(v, v^*) - f(\rho^{-1}r, r_0) \} \geq 0 \right) \leq \frac{\rho}{1-\rho} e^{-x}.$$

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

$$f(r, r_0) = 3\nu_0 r \sqrt{2x + 4p + 2 \log(r/r_0)}.$$

Idea: split  $\mathcal{B}_{\mathbf{r}^*}(\mathbf{v}^*)$  into slices  $\mathcal{B}_{\mathbf{r}_k}(\mathbf{v}^*) \setminus \mathcal{B}_{\mathbf{r}_{k-1}}(\mathbf{v}^*)$  and apply Theorem 9 to each slice.

By (3) and Theorem 9 for any  $\mathbf{r} > \mathbf{r}_0$

$$\begin{aligned} & I\!P\left( \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}^*) \setminus \mathcal{B}_{\rho\mathbf{r}}(\mathbf{v}^*)} \{\mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\mathbf{r}, \mathbf{r}_0)\} \geq 0 \right) \\ & \leq I\!P\left( \frac{1}{3\nu_0\mathbf{r}} \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}}(\mathbf{v}^*)} \mathcal{U}(\mathbf{v}, \mathbf{v}^*) \geq z_{\mathbb{H}}(\mathbf{x} + \log(\mathbf{r}/\mathbf{r}_0)) \right) \leq \frac{\mathbf{r}_0}{\mathbf{r}} e^{-\mathbf{x}}. \end{aligned} \quad (4)$$

Define  $\mathbf{r}_k = \mathbf{r}_0\rho^{-k}$  for  $k = 0, 1, 2, \dots$  and  $k^* \stackrel{\text{def}}{=} \log(\mathbf{r}^*/\mathbf{r}_0) + 1$ . By (4)

$$\begin{aligned} & I\!P\left( \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}^*}(\mathbf{v}^*) \setminus \mathcal{B}_{\mathbf{r}_0}(\mathbf{v}^*)} \left\{ \mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\rho^{-1}d(\mathbf{v}, \mathbf{v}^*), \mathbf{r}_0) \right\} \geq 0 \right) \\ & \leq \sum_{k=1}^{k^*} I\!P\left( \frac{1}{\mathbf{r}_k} \sup_{\mathbf{v} \in \mathcal{B}_{\mathbf{r}_k}(\mathbf{v}^*) \setminus \mathcal{B}_{\mathbf{r}_{k-1}}(\mathbf{v}^*)} \left\{ \mathcal{U}(\mathbf{v}, \mathbf{v}^*) - f(\mathbf{r}_k, \mathbf{r}_0) \right\} \geq 0 \right) \\ & \leq e^{-\mathbf{x}} \sum_{k=1}^{k^*} \rho^k \leq \frac{\rho}{1-\rho} e^{-\mathbf{x}}. \end{aligned}$$

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

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

$$\sup_{\mathbf{v} \in \mathcal{Y}} \log \mathbb{E} \exp \left\{ \lambda \boldsymbol{\gamma}^\top \nabla \mathcal{Y}(\mathbf{v}) \boldsymbol{\alpha} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad \lambda^2 \leq 2g^2. \quad (5)$$

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

$$\mathcal{Y}_o(r) = \{ \mathbf{v} \in \mathcal{Y}: \|\mathbf{v} - \mathbf{v}^*\| \leq r \}.$$

Condition (5) implies for any  $\mathbf{v} \in \mathcal{Y}_o(r)$  with  $\|\mathbf{v}\| \leq r$  and  $\|\boldsymbol{\gamma}\| = 1$  in view of  $\mathcal{Y}(\mathbf{v}^*) = 0$

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \boldsymbol{\gamma}^\top \mathcal{Y}(\mathbf{v}) \right\} \leq \frac{\nu_0^2 \lambda^2 \|\mathbf{v}\|^2}{2r^2}, \quad \lambda^2 \leq 2g^2; \quad (6)$$

Use the representation

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

This implies for  $\mathcal{Y}_o(r) = \{\mathbf{v} \in \mathcal{Y} : \|\mathbf{v} - \mathbf{v}^*\| \leq r\}$

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

Consider a bivariate process  $\mathbf{u}^\top \gamma(\mathbf{v})$  of  $\mathbf{u} \in \mathbb{R}^q$  and  $\mathbf{v} \in \Upsilon \subset \mathbb{R}^p$ .

By definition  $\mathbb{E}\mathbf{u}^\top \gamma(\mathbf{v}) = 0$ . Further, for  $\gamma = \mathbf{u}/\|\mathbf{u}\|$

$$\nabla_{\mathbf{u}} [\mathbf{u}^\top \gamma(\mathbf{v})] = \gamma(\mathbf{v}), \quad \nabla_{\mathbf{v}} [\mathbf{u}^\top \gamma(\mathbf{v})] = \mathbf{u}^\top \nabla \gamma(\mathbf{v}) = \|\mathbf{u}\| \gamma^\top \nabla \gamma(\mathbf{v})$$

Suppose that  $\mathbf{u} \in \mathbb{R}^q$  and  $\mathbf{v} \in \Upsilon$  are such that  $\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 \leq 2r^2$ . By the Hölder inequality, (6), and (5), it holds for  $\|\gamma\| = \|\alpha\| = 1$  and  $\mathbf{v} \in \Upsilon_r(r)$

$$\begin{aligned} & \log \mathbb{E} \exp \left\{ \frac{\lambda}{2r} (\gamma, \alpha)^\top \nabla [\mathbf{u}^\top \gamma(\mathbf{v})] \right\} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \gamma^\top \gamma(\mathbf{v}) \right\} + \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \mathbf{u}^\top \nabla \gamma(\mathbf{v}) \alpha \right\} \\ & \leq \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \gamma^\top \gamma(\mathbf{v}) \right\} + \frac{1}{2} \log \mathbb{E} \exp \left\{ \frac{\lambda}{r} \|\mathbf{u}\| \gamma^\top \nabla \gamma(\mathbf{v}) \alpha \right\} \\ & \leq \frac{\nu_0^2 \lambda^2}{4r^2} (\|\mathbf{v}\|^2 + \|\mathbf{u}\|^2) \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g. \end{aligned}$$

### Theorem

Let a random  $p$ -vector process  $\mathcal{Y}(\mathbf{v})$  for  $\mathbf{v} \in \Upsilon \subseteq \mathbb{R}^p$  fulfill  $\mathcal{Y}(\mathbf{v}^*) = 0$ ,  $E\mathcal{Y}(\mathbf{v}) \equiv 0$ , and the condition (5) be satisfied. Then for each  $r$  and any  $x \geq 1/2$ , it holds

$$P\left\{\sup_{v \in \Upsilon_r} \|\mathcal{Y}(v)\| > 6\nu_0 r z_{\mathbb{H}}(x)\right\} \leq e^{-x},$$

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

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

### 1 Introduction. Fisher and Wilks expansions

- Fisher and Wilks expansions
- The case of a linear model
- Expansions vs asymptotic results

### 2 Fisher and Wilks: Main steps

- Local quadraticity of  $\mathcal{EL}(\theta)$
- Local linear approximation of the stochastic term
- Local linear approximation of the gradient and the “Fisher” trick
- Local quadratic approximation of the log-likelihood and the “Wilks” trick
- Concentration and large deviation for  $\tilde{\theta}$
- A sharp bound for  $\|\xi\|^2$
- An upper function for the stochastic component

### 3 Examples

- Summary
- An i.i.d. case
- Generalized linear models
- Linear median (quantile) regression
- Conditional Moment Restriction (CMR)

$$\tilde{\boldsymbol{\theta}} \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta}), \quad \boldsymbol{\theta}^* \stackrel{\text{def}}{=} \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \mathbb{E} L(\boldsymbol{\theta})$$

### Theorem

On a set  $\Omega(x)$  with  $\mathbb{P}(\Omega(x)) \geq 1 - Ce^{-x}$

$$\|D_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) - \boldsymbol{\xi}\| \leq \diamond(x),$$

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

with

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

Here  $\diamond(x)$  and  $\Delta(x)$  are explicit error terms.

- Expansion of  $L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}^*)$

$$\begin{aligned} L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &= \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*) + \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla \zeta(\boldsymbol{\theta}^*). \\ &= \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\theta}^*) + \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top. \end{aligned}$$

- Taylor of the second order for  $\mathbb{E}L(\boldsymbol{\theta})$  around  $\boldsymbol{\theta}^*$ ;
- Local constant approximation of  $\nabla \zeta(\boldsymbol{\theta})$  – empirical processes theory for a vector stochastic process; involve the entropy function  $z_{\mathbb{H}}(x) \asymp \sqrt{p+x}$
- a sharp bound for the squared norm  $\|\boldsymbol{\xi}\|^2$ ;  $\mathbb{P}(\|\boldsymbol{\xi}\| \geq z(B, x)) \leq 2e^{-x}$  and  $z(B, x) \asymp \sqrt{\text{tr}(B) + x}$  for the “sandwich” matrix  $B$ ;
- “upper function” device for a centered stochastic process (remainder)  $\zeta(\boldsymbol{\theta}, \boldsymbol{\theta}^*) - \boldsymbol{\xi}^\top D_0(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top$  on an unbounded set  $\Theta$ . Disappears for linear models.

## Main steps

---

- Concentration and large deviations: identify  $r_0$  s.t.

$$\mathbb{P}(\tilde{\boldsymbol{\theta}} \notin \Theta_0(r_0)) \leq e^{-x},$$

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

- Local quadratic approximation of the expected log-likelihood:

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

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

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

- Overall error of the Fisher expansion  $r_0\{\delta(r_0) + \varrho(r_0, x)\}$ ,  
of the Wilks  $r_0^2\{\delta(r_0) + \varrho(r_0, x)\}$ .

## Conditions

(**L<sub>0</sub>**) For each  $r \leq r_0$ , there is a constant  $\delta(r) \leq 1/2$  such that it holds for any  $\theta \in \Theta_0(r) = \{\theta \in \Theta : \|D_0(\theta - \theta^*)\| \leq r\}$ :

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

(**ED<sub>2</sub>**) There exist a value  $\omega > 0$  and for each  $r > 0$ , a constant  $g(r) > 0$  such that  $\zeta(\theta) \stackrel{\text{def}}{=} L(\theta) - \mathbb{E}L(\theta)$  satisfies for any  $\theta \in \Theta_0(r)$ :

$$\sup_{\gamma_1, \gamma_2 \in \mathbb{R}^p} \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \frac{\gamma_1^\top \nabla^2 \zeta(\theta) \gamma_2}{\|D_0 \gamma_1\| \cdot \|D_0 \gamma_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g(r).$$

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

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

(**L**) For each  $r$ , there exists  $b(r) > 0$  such that  $r b(r) \rightarrow \infty$  as  $r \rightarrow \infty$  and

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

Consider  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ .

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

Yields

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

$$\tilde{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} L(\theta),$$

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E} L(\theta).$$

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

$$D_n^2 = n \mathbb{F}_{\theta^*}.$$

(for simplicity  $p = 1$ )

► Smoothness:

- $\nabla^2 \mathbb{E}\ell(Y_1, \boldsymbol{\theta})$  Lipschitz continuous in  $\boldsymbol{\theta}$ ;
- $\mathbb{E} \exp\{\lambda_0 \ell'(Y_1, \boldsymbol{\theta})\} \leq C$
- $\mathbb{E} \exp\{\lambda_0 \ell''(Y_1, \boldsymbol{\theta})\} \leq C$

► Identifiability:

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

Then the conditions are fulfilled with  $g^2 \approx n\lambda_0$  and  $b(r) \geq b_0 > 0$ .

## Checking ( $ED_0$ )

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

Let

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

and

$$\log \exp\{\lambda \mathbf{v}_0^{-1} \nabla \zeta_i(\boldsymbol{\theta}^*)\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad |\lambda| \leq g_0$$

Then for  $|\lambda| \leq g_0 n^{1/2}$

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

Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top \sim \mathbb{I}\mathcal{P}$ , a sample of independent r.v.s.

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

- $\Psi_i$ , given factors in  $\mathbb{R}^p$ ,
- $(P_{\mathbf{v}})$ , an exponential family with canonical parametrization,  $\ell(y, \mathbf{v}) = y\mathbf{v} - d(\mathbf{v})$ ,
- $\boldsymbol{\theta} \in \mathbb{R}^p$ , unknown parameter.

MLE:

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

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

with  $f_i = \mathbb{E} Y_i$ .

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

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

$$\zeta(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) - \mathbb{E}L(\boldsymbol{\theta}) = \left( \sum_{i=1}^n \varepsilon_i \Psi_i \right)^\top \boldsymbol{\theta}$$

$\nabla^2 \zeta(\boldsymbol{\theta}) \equiv 0$  and ( $ED_2$ ) automatically;

The Fisher information  $D_0^2$  depends on  $\mathbb{P}$  only through  $\boldsymbol{\theta}^*$ :

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

The same for the vector  $\xi$ :

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

Sufficient conditions:

- $d''(\Psi_i^\top \theta)$  uniformly continuous in  $\theta$  over  $i = 1, \dots, n$ ; here  $\ell(y, v) = yv - d(v)$ ;
- for some fixed matrices  $v_i^2$  and  $\lambda_0 > 0$

$$\mathbb{E} \exp\{\lambda_0 v_i^{-1} \varepsilon_i\} \leq c$$

- the matrix  $V_0^2 = \sum_i v_i^2$  fulfills

$$V_0^2 \leq a^2 D_0^2.$$

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

## Median (quantile) regression

Consider a median linear regression

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

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

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

MLE = LAD

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

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

## Checking the conditions

---

Sufficient conditions:

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

$$\sup_i \sup_{|u| \leq t} \left| \frac{f_i(u)}{f_i(0)} - 1 \right| \leq \omega(t), \quad \text{small for } t \text{ small};$$

the sample size  $n$  satisfies

$$n \geq C p$$

Challenges:

- $\nabla L(\boldsymbol{\theta})$  exists but **discontinuous**;
- $\nabla^2 L(\boldsymbol{\theta})$  is a delta-function;
- $I\!E L(\boldsymbol{\theta})$  is **Lipschitz** in  $\boldsymbol{\theta}$ ; we need that  $I\!E L(\boldsymbol{\theta}, \boldsymbol{\theta}^*)$  grows **faster than linear**.

Observed  $Z_i = (X_i, Y_i)$ .

Conditional estimating equations (or moment restrictions)

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

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:  $g(Z, \boldsymbol{\theta}) = \mathbb{1}\{Y - f(X, \boldsymbol{\theta}) \leq 0\} - \tau$  for a quantile of order  $\tau$ ;
3. linear transformation regression models:  $g(Z, \boldsymbol{\theta}) = h(Y, \boldsymbol{\eta}) - X^\top \boldsymbol{\beta}$  and  $\boldsymbol{\theta} = (\boldsymbol{\eta}^\top, \boldsymbol{\beta}^\top)^\top$ ;
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**:

$$\mathbb{E}[A(X)g(Z, \theta_0)] = 0 \quad \text{a.s.}$$

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):

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

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**.

## Approach

---

Let  $Z$  be the observed data. Define

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

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

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

## Correction for the noise energy

Simple calculus yields the expectation

$$\mathbb{E}M(\boldsymbol{\theta}) = \sum_{i,j} b_i(\boldsymbol{\theta})b_j(\boldsymbol{\theta})w_{ij} + \sum_i \mathbb{E}\varepsilon_i^2(\boldsymbol{\theta}) w_{ii}, \quad (7)$$

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 (7):

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

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

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

Define also

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

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

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

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

## Approach

---

Represent

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

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

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

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

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

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

## A bound for $\|A\{\varepsilon(\theta) - \varepsilon(\theta^*)\}\|$

### Theorem

Suppose that for any  $\theta \in \Theta_0(x_0)$  and each  $i = 1, \dots, n$  and any unit vector  $\gamma \in I\!\!R^p$

$$\log I\!\!E \exp \left\{ \lambda \gamma^\top \nabla \varepsilon_i(\theta) \right\} \leq \nu_0^2 \lambda^2 / 2, \quad \lambda^2 \leq 2g^2,$$

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

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

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

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

Here  $\mathbb{H}_1 = 2\mathbb{H}_1(A)$  and  $\mathbb{H}_2 = \mathbb{H}_2(A) + 2c_1 p$  with

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

## A leading term

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

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

It is obvious that

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

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

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

Under PA, it holds  $b_i(\boldsymbol{\theta}^*) \equiv 0$  and

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

$$g_i(\boldsymbol{\theta}) = b_i(\boldsymbol{\theta}) + \varepsilon_i(\boldsymbol{\theta}),$$

$$\boldsymbol{\theta}^* = \operatorname{arginf}_{\boldsymbol{\theta}} \mathbf{b}(\boldsymbol{\theta})^\top W \mathbf{b}(\boldsymbol{\theta}),$$

$$\tilde{\boldsymbol{\theta}} = \operatorname{arginf}_{\boldsymbol{\theta}} \left\{ \mathbf{g}(\boldsymbol{\theta})^\top W \mathbf{g}(\boldsymbol{\theta}) - S(\boldsymbol{\theta}) \right\}.$$

### Theorem

Suppose that for any  $\boldsymbol{\theta} \in \Theta_0(x_0)$  and each  $i = 1, \dots, n$  and any unit vector  $\gamma \in \mathbb{R}^p$

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

the functions  $b_i(\boldsymbol{\theta})$  are twice continuously differentiable uniformly in  $i$ , and the identifiability condition holds.

Then  $\tilde{\boldsymbol{\theta}}$  is root-n consistent and semi parametrically efficient estimator of  $\boldsymbol{\theta}^*$ .