# Outline
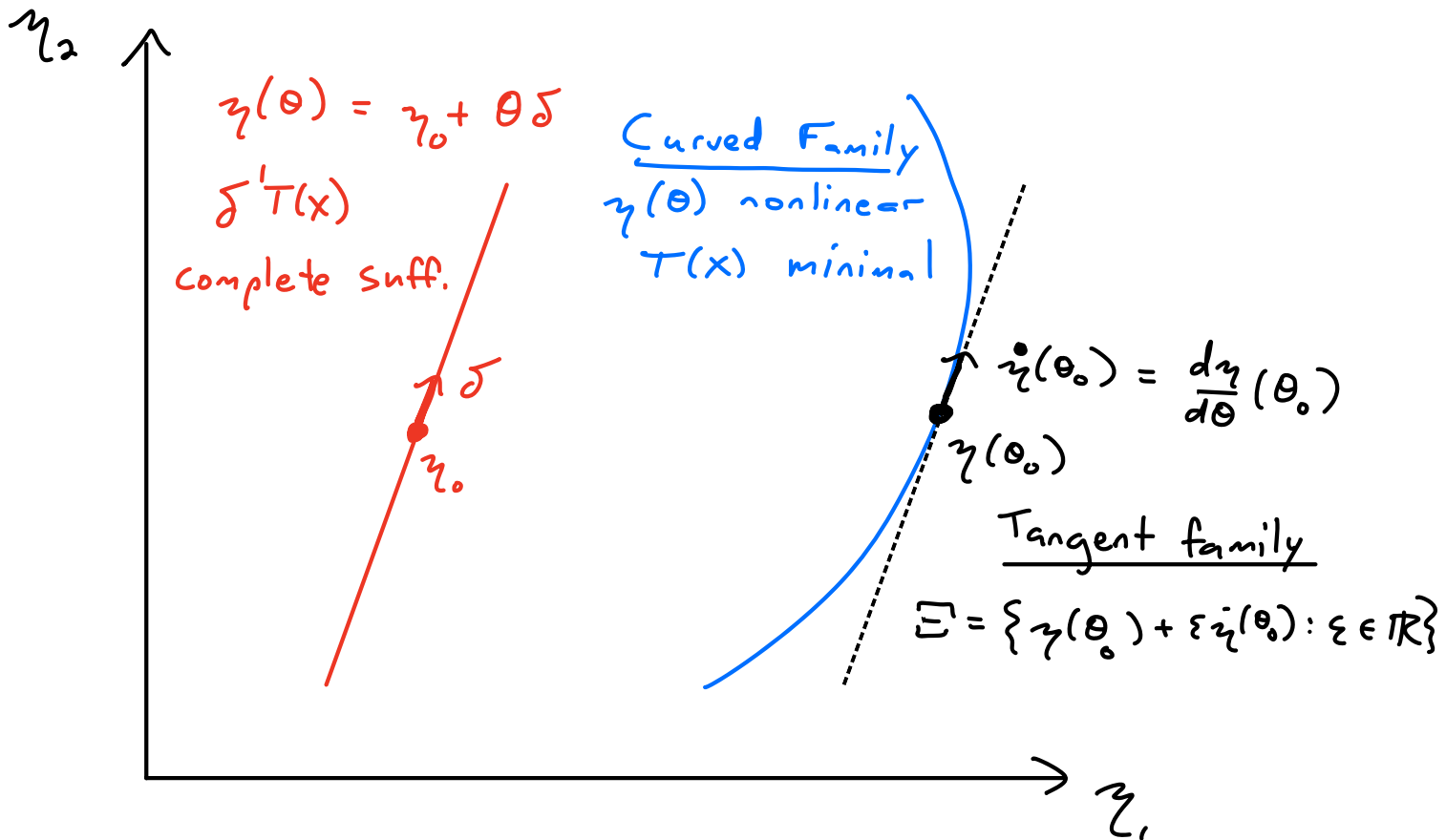
1) Score function

2) Fisher information

3) Cramér-Rao Lower Bound

4) Examples

# Motivation : Tangent family

$$p_\theta(x) = e^{\eta(\theta)' T(x) - A(\eta(\theta))} h(x) \qquad \eta: \mathbb{R} \to \mathbb{R}^2$$

$$\Xi = \{\eta(\theta) : \theta \in \mathbb{R}\}$$

$\eta_2$ (axis)

$\eta(\theta) = \eta_0 + \theta \delta$

$\delta' T(x)$

complete suff.

Curved Family
$\eta(\theta)$ nonlinear
$T(x)$ minimal

$\delta$

$\eta_0$

$\dot{\eta}(\theta_0) = \frac{d\eta}{d\theta}(\theta_0)$

$\eta(\theta_0)$

Tangent family

$$\Xi = \{\eta(\theta_0) + \varepsilon \dot{\eta}(\theta_0) : \varepsilon \in \mathbb{R}\}$$

$\eta_1$ (axis)

$$q_\varepsilon(x) = e^{(\eta(\theta_0) + \varepsilon \dot{\eta}(\theta_0))' T(x) - A(\cdots)} h(x)$$

$$= e^{\underbrace{\varepsilon \dot{\eta}(\theta_0)'(T(x) - \mathbb{E}_{\theta_0} T)}_{S_{\theta_0}(x)} - B(\varepsilon)} k(x)$$

Complete sufficient for tangent family at $\theta_0$
Called  <u>Score function</u>

# Score function

Assume $\mathcal{P}$ has densities $\rho_\theta$ wrt $\mu$, $\Theta \subseteq \mathbb{R}^d$

Common support: $\{x : \rho_\theta(x) > 0\}$ same $\forall \theta$

Recall $\ell(\theta; x) = \log \rho_\theta(x)$,

Thought of as <u>random function of $\theta$</u>

<u>Def</u> The <u>score</u> is $\nabla \ell(\theta; x)$; plays a key role in many areas of statistics, esp. asymptotics.

Can think of as "local complete sufficient statistic":

$$\rho_{\theta_0 + z}(x) = e^{\ell(\theta_0 + z; x)}$$

$$\approx e^{z' \nabla \ell(\theta_0; x)} \rho_{\theta_0}(x) \qquad \text{for } z \approx 0$$

<u>Differential identities</u>: (assuming enough regularity)

$$1 = \int_{\mathcal{X}} e^{\ell(\theta; x)} d\mu(x)$$

$$\frac{\partial}{\partial \theta_j} \Rightarrow \quad 0 = \int \frac{\partial}{\partial \theta_j} \ell(\theta; x)\, e^{\ell(\theta; x)} d\mu(x)$$

$$\Rightarrow \quad \mathbb{E}_\theta \left[ \nabla \ell(\theta; x) \right] = 0$$

only true if these are the same value of $\theta$!

$$\frac{\partial}{\partial \theta_k} \Rightarrow \quad 0 = \int \left( \frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} + \frac{\partial \ell}{\partial \theta_j} \cdot \frac{\partial \ell}{\partial \theta_k} \right) e^{\ell} \, d\mu$$

$$= \mathbb{E}_{\theta}\left[ \frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} \right] + \mathbb{E}_{\theta}\left[ \frac{\partial \ell}{\partial \theta_j} \quad \frac{\partial \ell}{\partial \theta_k} \right]$$

$$\Longrightarrow \quad J(\theta) = \mathrm{Var}_{\theta}\left[ \nabla \ell(\theta; X) \right] = \mathbb{E}_{\theta}\left[ -\nabla^2 \ell(\theta; X) \right]$$

<span style="color:red">↖ ↗ same $\theta$</span>  <span style="color:red">↖ ↗ same $\theta$</span>

Called "Fisher Information"

<span style="color:red">[</span> It is possible to extend this definition to certain cases where $\ell$ is not even differentiable, e.g. Laplace location family, but for our purposes we can just assume "sufficient regularity." <span style="color:red">]</span>

Try with another statistic $\delta(X)$, let

$$g(\theta) = \mathbb{E}_{\theta}\left[ \delta(X) \right] \qquad (\text{"unbiased estimator"})$$

$$g(\theta) = \int \delta e^{\ell} d\mu$$

$$\Rightarrow \quad \nabla g(\theta) = \int \delta \nabla \ell \, e^{\ell} d\mu = \mathbb{E}_{\theta}\left[ \delta(X) \nabla \ell(\theta; X) \right]$$

$$\overset{\underset{\text{<span style="color:red">Since $\mathbb{E}\nabla \ell = 0$</span>}}{\nearrow}}{=} \mathrm{Cov}_{\theta}\left( \delta(X), \nabla \ell(\theta; X) \right)$$

Combining these results with Cauchy-Schwarz gives us the $\underline{\text{Cramér-Rao Lower Bound}}$ or $\underline{\text{Information Lower Bound}}$:

$\underline{\text{1-param}}$: $\quad Var_\theta(\delta) \cdot Var_\theta(\dot{\ell}(\theta;x)) \geq Cov_\theta(\delta, \dot{\ell}(\theta;x))^2$

$$\Rightarrow \quad Var_\theta(\delta) \gtrsim \dot{g}(\theta)^2 / J(\theta)$$

$\underline{\text{Multivariate}}$: $\quad \theta \in \mathbb{R}^d, \quad g(\theta), \delta(x) \in \mathbb{R}$

$$Var_\theta(\delta) \geq \nabla g(\theta)' J(\theta)^{-1} \nabla g(\theta)$$

$\underline{\text{Proof}}$:

$$Var_\theta(\delta) \cdot a' J(\theta) a = Var_\theta(\delta) \, Var(a' \nabla \ell(\theta))$$
$$\geq Cov_\theta(\delta, a' \nabla \ell(\theta))^2$$
$$= a' \nabla g \, \nabla g' a \quad, \text{ for all } \quad a \in \mathbb{R}^d$$

$$\Rightarrow \quad Var_\theta(\delta) \geq \max_{a \neq 0} \frac{a' \nabla g \nabla g' a}{a' J(\theta) a} \underset{\text{Exercise}}{\boxed{=}} \nabla g' J(\theta)^{-1} \nabla g$$

$u = J(\theta)^{1/2} a$

$\max_u \frac{u' J^{-1/2} \nabla g \nabla g' J^{-1/2} u}{u' u}$

$u = J^{-1/2} \nabla g \qquad a = J^{-1} \nabla g$

$\underline{\text{Interp}}$: If $g(\theta)$ is estimand, no unbiased estimator can have smaller variance than $\nabla g(\theta)' J(\theta)^{-1} \nabla g(\theta)$

**Ex.** : (i.i.d. sample)

$$X_1, \ldots, X_n \overset{iid}{\sim} p_\theta^{(1)}(x) \qquad \theta \in \Theta \subseteq \mathbb{R}^d$$

$p_\theta$ "regular": common support, finite derivative wrt $\theta$

$$X \sim p_\theta(x) = \prod_i p_\theta^{(1)}(x_i)$$

Let $\ell_1(\theta; x_i) = \log p_\theta^{(1)}(x_i)$

$$\ell(\theta; x) = \sum_i \ell_1(\theta; x_i)$$

$$J(\theta) = \text{Var}_\theta\left(\nabla \ell(\theta; X)\right)$$

$$= \text{Var}_\theta\left(\sum_i \nabla \ell_1(\theta; X_i)\right)$$

$$= n J_1(\theta) \qquad \text{where } J_1(\theta) \text{ is Fisher info in single observation}$$

$\Rightarrow$ Lower bound scales like $n^{-1}$ (SD $\approx n^{-1/2}$ for "regular" families)

# Efficiency

CRLB is not nec. attainable.

We define the efficiency of an unbiased estimator as:

$$\text{eff}_\theta(\delta) = \frac{CRLB}{\text{Var}_\theta(\delta)} \left( = \frac{1/J(\theta)}{\text{Var}_\theta(\delta)} \quad \text{if } g(\theta) = \theta \in \mathbb{R} \right)$$

$$\text{eff}_\theta(\delta) \leq 1$$

We say $\delta(X)$ is <u>efficient</u> if $\text{eff}_\theta(\delta) = 1 \quad \forall \theta$

Depends on $\text{Corr}_\theta(\delta(X), \nabla \ell(\theta; X))$ :

$$\text{eff}_\theta(\delta) = \frac{\text{Cov}_\theta^2(\delta(X), \dot{\ell}(\theta; X))}{\text{Var}_\theta(\delta) \cdot \text{Var}_\theta(\dot{\ell}(\theta))}$$

$$= \text{Corr}_\theta^2(\delta, \dot{\ell}(\theta))$$

$$\leq 1$$

$\delta(X)$ is efficient $\iff \text{Corr}_\theta^2(\delta, \dot{\ell}(\theta)) = 1 \quad \forall \theta$

Rarely achieved in finite samples but we can approach it asymptotically as $n \to \infty$

**Ex.** Exponential Families

$$p_\eta(x) = e^{\eta' T(x) - A(\eta)} h(x)$$

$$\ell(\eta; x) = \eta' T(x) - A(\eta) + \log h(x)$$

$$\nabla \ell(\eta; x) = T(x) - \nabla A(\eta)$$

$$= T(x) - \mathbb{E}_\eta T(x)$$

$$\text{Var}_\eta(\nabla \ell(\eta)) = \text{Var}_\eta(T(x)) = \nabla^2 A(\eta)$$

$$\nabla^2 \ell(\eta; x) = - \nabla^2 A(\eta)$$

$$\mathbb{E}_\eta\left[-\nabla^2 \ell(\eta; x)\right] = \nabla^2 A(\eta) \qquad \checkmark$$

So any unbiased est. of $\eta$ has

$$\text{Var}_\eta(\delta) \geq \nabla^2 A(\eta)^{-1}$$

Curved family: $\quad p_\theta(x) = e^{\gamma(\theta)' T(x) - B(\theta)} h(x), \quad \theta \in \mathbb{R}$
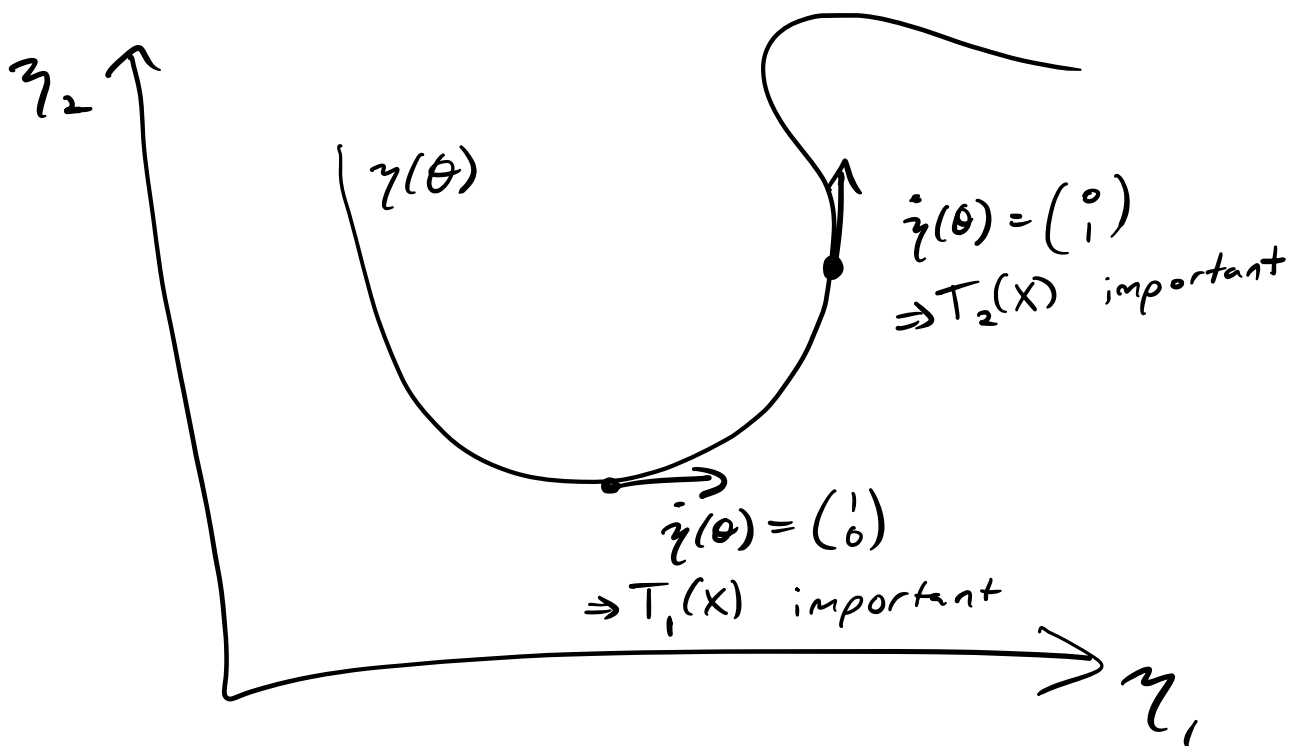
$$B(\theta) = A(\gamma(\theta))$$

$$\ell(\theta; x) = \gamma(\theta)' T(x) - B(\theta) + \log h(x)$$

$$\dot\ell(\theta; X) = \dot\gamma(\theta)' T(x) - \dot\gamma(\theta)' \nabla_\gamma A(\gamma(\theta))$$

$$= \dot\gamma(\theta)' \left( T(x) - \nabla_\gamma A(\gamma(\theta)) \right)$$

$$= \dot\gamma(\theta)' \left( T(X) - \mathbb{E}_\theta T(x) \right)$$

$$\implies \dot\gamma(\theta)' T(X) \quad \text{is} \quad \text{``locally complete suff. stat.''}$$



$\eta_2$

$\gamma(\theta)$

$\dot\gamma(\theta) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$
$\Rightarrow T_2(X)$ important

$\dot\gamma(\theta) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$\Rightarrow T_1(X)$ important

$\eta_1$

# Fisher info as local metric

Kullback-Leibler Divergence

$$D_{KL}(p \| q) = \mathbb{E}_p \left[ \log p(X) - \log q(X) \right]$$

$$= \int \log \left( \frac{p}{q} \right) p \, d\mu$$

Distance between two distributions

Parametric model

$$D_{KL}(\theta^* \| \theta) = D_{KL}\left( p_{\theta^*} \| p_\theta \right)$$

$$= \int (l(\theta^*) - l(\theta)) e^{l(\theta^*)} \, d\mu$$

Standard "distance" between two distributions

$\theta^*$ "real" distribution, function of $\theta$

Maximized at $\theta = \theta^*$ :

$$\frac{\partial}{\partial \theta_j} D_{KL}(\theta^* \| \theta) = -\int \frac{\partial \ell}{\partial \theta_j}(\theta) \, e^{\ell(\theta^*)} \, d\mu$$

$$= 0 \quad \text{at } \theta = \theta^*$$

$$\frac{\partial^2}{\partial \theta_j \partial \theta_k} D_{KL}(\theta^* \| \theta) = -\int \frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k}(\theta) \, e^{\ell(\theta^*)} \, d\mu$$

$$= + J(\theta^*)_{jk} \quad \text{at } \theta = \theta^*$$

$d = 1$ :