

NFST 434 - Ex. session 6. M-estimators

$X_1, \dots, X_n$  random sample from  $F$  with unknown parameter  $\theta \in \Theta \subset \mathbb{R}^p$ .

$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \rho(X_i, \theta)$  M-estimator

$O_p = \sum_{i=1}^n \psi(X_i, \hat{\theta})$  Z-estimator (often  $\psi(x, \theta) = \frac{\partial \rho(x, \theta)}{\partial \theta}$ )

$\hat{\theta}$  estimator  $\theta_x = \underset{\theta}{\operatorname{argmin}} E \rho(X_1, \theta)$  and  $O_p = E \psi(X_1, \theta_x)$ , respectively.

Under [25]-[26],  $\Gamma(\theta) = E D_{\psi}(X_1, \theta)$ ,  $\Sigma(\theta_x) = E(\psi(X_1, \theta_x) \psi(X_1, \theta_x)^T)$  from Theorem 9

$\Gamma_n(\hat{\theta} - \theta_x) = -\frac{1}{\Gamma_n} \sum_{i=1}^n \Gamma^{-1}(\theta_x) \psi(X_i, \theta_x) + o_p(1)$ , i.e.  $\Gamma_n(\hat{\theta} - \theta_x) \xrightarrow{d} N_p(O_p, \Gamma^{-1}(\theta_x) \Sigma(\theta_x) \Gamma^{-1}(\theta_x))$ .

Ex.  $N(\mu, \sigma^2)$  MLE  $L(\mu) = c \cdot \sigma^{-n} \exp(-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2)$ ,  $\ell(\mu) = c' - \frac{1}{2} \sum (x_i - \mu)^2$ ,  $\rho(x, \theta) = (x - \mu)^2$

i.e.  $\psi(x, \theta) = \frac{\partial}{\partial \mu} (x - \mu)^2 = 2(x - \mu)$  and  $\bar{X}$  is an M-estimator of  $E(X_1 - \mu) = 0$ ,  $\theta_x = EX = \mu$ .

now if the true distribution is not  $N(\mu, \sigma^2)$ , we get  $E D_{\psi}(X_1, \theta) = -1 = \Gamma(\mu)$ ,

$\operatorname{var} \psi(X_1, \theta_x) = \operatorname{var}(X_1 - \mu_x) = \operatorname{var} X_1$  and  $\Gamma_n(\bar{X} - \mu) \xrightarrow{d} N(0, \operatorname{var} X_1)$  no matter

what the true distribution of  $X_1$  is. (CLV.) By Example 34 this choice of  $\theta_x = EX$  minimizes the Kullback-Leibler divergence of the family  $\{N(\mu, \sigma^2), \mu \in \mathbb{R}\}$  from the true distribution  $F$  of  $X_1$ , i.e.  $\theta_x = \underset{\theta}{\operatorname{argmin}} \int_{\mathbb{R}} \log \left[ \frac{f(x)}{f(x, \theta)} \right] \cdot f(x) d\mu(x)$ . This holds generally for M-estimators

given as MLE in systems of densities.  $L(\theta)$  likelihood,  $\rho(x, \theta) = -\ell(x, \theta)$ ,  $\psi(x, \theta) = \frac{\partial \ell(x, \theta)}{\partial \theta}$ .

Ex.  $N(\mu, \sigma^2)$   $\ell(\mu, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$ ,  $\theta = (\mu, \sigma^2)^T$ ,  $\rho(x, \theta) = \frac{(x - \mu)^2}{\sigma^2} + \log \sigma^2$

$\psi(x, \theta) = \begin{pmatrix} x - \mu \\ (x - \mu)^2 - \sigma^2 \end{pmatrix}$  solves  $\hat{\mu} = \bar{X}$   
 $\hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{X})^2$ , population  $E(X - \mu) = 0 \Rightarrow \mu_x = EX$   
 $E(X - \mu)^2 = \sigma^2 \Rightarrow \sigma_x^2 = \operatorname{var} X$

$E D_{\psi}(X_1, \theta_x) = E \begin{pmatrix} -1 & 0 \\ -2(X_1 - \mu_x) & -1 \end{pmatrix} = \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} = \Gamma(\theta_x)$   
 $\Sigma(\theta_x) = \begin{pmatrix} \sigma_x^2 & E(X - \mu_x)^3 \\ E(X - \mu_x)^3 & E(X - \mu_x)^4 - \sigma^4 \end{pmatrix}$   
 $E(X - \mu_x)(X - \mu_x)^2 - \sigma_x^3 = E(X - \mu_x)^3$   
 $E[(X - \mu_x)^2 - \sigma^2]^2 = E(X - \mu_x)^4 - \sigma^4$

$\Gamma_n \left( \begin{pmatrix} \bar{X} \\ \hat{\sigma}^2 \end{pmatrix} - \begin{pmatrix} \mu_x \\ \sigma_x^2 \end{pmatrix} \right) \xrightarrow{d} N_2 \left( O_2, \begin{pmatrix} \operatorname{var} X_1 & E(X_1 - \mu_x)^3 \\ E(X_1 - \mu_x)^3 & E(X_1 - \mu_x)^4 - \sigma^4 \end{pmatrix} \right)$  again CLV. (\*)  
 $= N_2 \left( O_2, \begin{pmatrix} \operatorname{var} X_1 & 0 \\ 0 & 2\sigma^4 \end{pmatrix} \right)$  under normality.

MLE: for the choice  $\psi_2(x, \theta) = \begin{pmatrix} \frac{x - \mu}{\sigma^2} \\ \frac{(x - \mu)^2}{2\sigma^4} - \frac{1}{2\sigma^2} \end{pmatrix}$  we get the same solution  $\hat{\mu} = \bar{X}$   
 $\hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{X})^2$

but  $\Gamma(\theta_x) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}$ ,  $\Sigma(\theta_x) = \begin{pmatrix} \frac{1}{\sigma^2} & \frac{E(X_1 - \mu)^3}{2\sigma^3} \\ \frac{E(X_1 - \mu)^3}{2\sigma^3} & \frac{E(X_1 - \mu)^4 - \sigma^4}{4\sigma^8} \end{pmatrix}$  (under normality  $\Sigma(\theta_x) = \Gamma(\theta_x) = I(\theta_x)$ )  
 as follows from general theory

and one gets the sandwich estimator of the as. variance of  $(\bar{X}, \hat{\sigma}^2)^T$  Fisher's inform. matrix

$\Gamma^{-1}(\theta_x) \Sigma(\theta_x) \Gamma^{-1}(\theta_x)$  with  $\theta_x$  replaced by  $\hat{\theta}$ . This estimator holds true also

under misspecification and is the same as in (\*).



Ex: M-estimators and  $\Delta$ -theorem Influence also about  $\sigma$ . Parameter  $(\mu, \sigma^2, \theta)$  for  $\theta = \sqrt{\sigma^2}$ .

Estimating function  $\psi(\mu, \sigma^2, \theta) = \begin{pmatrix} x - \mu \\ (x - \mu)^2 - \sigma^2 \\ \sqrt{\sigma^2} - \theta \end{pmatrix}$   $\theta_x = \sigma_x$  estimating equations  $\Leftrightarrow$  MLE which is invariant to reparametrization

$$\Gamma(\mu, \sigma^2, \theta) = E \begin{pmatrix} -1 & 0 & 0 \\ -2(x-\mu) & -1 & 0 \\ 0 & \frac{1}{2\sqrt{\sigma^2}} & -1 \end{pmatrix} = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & \frac{1}{2\sigma} & -1 \end{pmatrix}, \quad \Gamma^{-1}(\mu, \sigma^2, \theta) = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & -1/2\sigma & -1 \end{pmatrix}$$

$$\Sigma(\mu, \sigma^2, \theta) = \text{var} \begin{pmatrix} x - \mu \\ (x - \mu)^2 - \sigma^2 \\ \sigma - \theta \end{pmatrix} = \begin{pmatrix} \sigma^2 & \mu_3 & 0 \\ \mu_3 & \mu_4 - \sigma^4 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

$$\mu_i = E(X - \mu)^i$$

$$\Gamma_m \left( \begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \\ \hat{\theta} \end{pmatrix} - \begin{pmatrix} \mu \\ \sigma^2 \\ \theta \end{pmatrix} \right) \xrightarrow{d} N_3 \left( 0, \Gamma^{-1}(\mu, \sigma^2, \theta) \Sigma(\mu, \sigma^2, \theta) \Gamma^{-1}(\mu, \sigma^2, \theta) \right)$$

See the Mathematica script  
 $\hat{\Sigma}_m$  as an M-estimator

$$\begin{pmatrix} \sigma^2 & \mu_3 & \mu_3/2\sigma \\ \mu_3 & \mu_4 - \sigma^4 & (\mu_4 - \sigma^4)/2\sigma \\ \mu_3/2\sigma & (\mu_4 - \sigma^4)/2\sigma & (\mu_4 - \sigma^4)/4\sigma^2 \end{pmatrix}$$

the same result as obtained using the  $\Delta$ -theorem

**Influence function.** considers a general estimator  $\hat{\theta}$  of  $\theta \in \mathbb{R}^p$ ,  $\theta = \theta(P)$  functional of the measure  $P$ . Let  $t \in [0, 1]$ ,  $x$  in the sample space and  $P_t := (1-t)P \oplus t\delta_x$  mixture of  $P$  and  $\delta_x$ . Influence function is of  $\theta(P)$  is the Gâteaux (directional) derivative of  $\theta(P)$  in direction  $\delta_x$ , i.e.  $IF(x) = \lim_{\varepsilon \rightarrow 0} \frac{\theta(P_\varepsilon) - \theta(P_0)}{\varepsilon} = \frac{\partial \theta_\varepsilon}{\partial \varepsilon} \Big|_{\varepsilon=0}$  for  $\theta_\varepsilon := \theta(P_\varepsilon)$

For  $\varepsilon$  small  $\theta(P_\varepsilon) \approx \theta(P_0) + \varepsilon \cdot IF(x)$ . quantifies the effect of an infinitesimal change of  $P$  in direction  $\delta_x$  on the (estimator)  $\theta(P)$ . Bounded IF implies resistance to outliers and measurement errors.  $\rightarrow$  Robustness

For M-estimator:  $\theta(P)$  solves  $E_P \psi(X, \theta(P)) = 0_P$ , i.e.  $E_{P_\varepsilon} \psi(X, \theta_\varepsilon) = 0 =$

$$= \int \psi(y, \theta_\varepsilon) dP_\varepsilon(y) = (1-\varepsilon) \int \psi(y, \theta_\varepsilon) dP(y) + \varepsilon \int \psi(y, \theta_\varepsilon) d\delta_x(y) = (1-\varepsilon) \int \psi(y, \theta_\varepsilon) dP(y) + \varepsilon \psi(x, \theta_\varepsilon)$$

$$\cdot \frac{\partial}{\partial \varepsilon} \text{ gives } 0 = - \int \underbrace{\psi(y, \theta_\varepsilon)}_{p \times 1} dP(y) + (1-\varepsilon) \underbrace{\int \frac{\partial \psi(y, \theta)}{\partial \theta^T}}_{p \times p} \Big|_{\theta_\varepsilon} dP(y) + \underbrace{\frac{\partial \theta_\varepsilon}{\partial \varepsilon}}_{p \times 1} + \psi(x, \theta_\varepsilon)$$

$$+ \varepsilon \cdot \frac{\partial}{\partial \theta^T} \psi(x, \theta) \Big|_{\theta_\varepsilon} \cdot \frac{\partial \theta_\varepsilon}{\partial \varepsilon} \quad \cdot \varepsilon \rightarrow 0+$$

$$\frac{\partial \theta_\varepsilon}{\partial \varepsilon} = \lim_{\varepsilon \rightarrow 0} \frac{\theta_{\varepsilon+\varepsilon} - \theta_\varepsilon}{\varepsilon} \quad \cdot \varepsilon \rightarrow 0+ \rightsquigarrow \frac{\partial \theta_\varepsilon}{\partial \varepsilon} \Big|_{\varepsilon=0} = IF(x)$$

under appropriate conditions  $\varepsilon \frac{\partial}{\partial \theta^T} \psi(x, \theta) \Big|_{\theta_\varepsilon} \frac{\partial \theta_\varepsilon}{\partial \varepsilon} \xrightarrow{\varepsilon \rightarrow 0} 0$

$$0 = - \int \underbrace{\psi(y, \theta_0)}_{0 \text{ from identification of } \theta_0, \text{ because } \theta_0 \text{ is given by } E_P \psi(X, \theta_0) = 0} dP(y) + \Gamma(\theta_0) \cdot IF(x) + \psi(x, \theta_0) + 0$$

$$IF(x) = - \Gamma^{-1}(\theta_0) \cdot \psi(x, \theta_0)$$

From Theorem 9 also  $\hat{\theta}_m - \theta_x = \frac{1}{m} \sum_{i=1}^m IF(X_i) + o_p \left( \frac{1}{\sqrt{m}} \right)$  and a bounded IF, or equivalently bounded  $\psi(x, \theta)$  in  $x$  results in resistant estimators.

Above  $\psi(x, \theta) = (x - \mu)^j$ ,  $j \geq 1 \rightsquigarrow$  unbounded IF,  $\bar{X}$  and  $\hat{\sigma}^2$  are affected by outliers



Ex. Robust location estimators

Ex. session 7

Cauchy distribution MLE  $f(x, \theta) = \frac{1}{\pi} \frac{1}{1+(x-\mu)^2}$   $l(\mu) = c - \log(1+(x-\mu)^2)$

$\psi(x, \mu) = -\frac{\partial l(\mu)}{\partial \mu} = \frac{2(x-\mu)}{1+(x-\mu)^2}$   $\mu_x = \operatorname{argmin}_{\mu} E \log(1+(X-\mu)^2)$

$\hat{\mu} = \operatorname{argmin}_{\mu} \frac{1}{n} \sum_{i=1}^n \log(1+(x_i-\mu)^2)$  computable only numerically, non-convex minimization  
if  $X \sim \text{Cauchy}(\mu)$

$\Gamma(\mu) = E D_{\psi}(X, \mu) = E \frac{2(1-(x-\mu)^2)}{(1+(x-\mu)^2)^2} = \frac{1}{2}$

$\Sigma(\mu_x) = \operatorname{var} \frac{2(X-\mu_x)}{1+(X-\mu_x)^2} = \frac{1}{2}$  if  $X \sim \text{Cauchy}$

if  $X \sim \text{Cauchy}$   $\Gamma_n(\hat{\mu} - \mu_x) \xrightarrow{d} N(0, 2)$

$|\psi(x, \mu)| \leq 1 \quad \forall x, \mu$  generalized to t-distributions

Laplace MLE  $f(x, \theta) = \frac{1}{2} \exp(-|x-\theta|)$

$l(\theta) = c - |x-\theta|$

$\theta_x = \operatorname{argmin}_{\theta} E |X-\theta|$   $\hat{\theta} = \operatorname{argmin}_{\theta} \frac{1}{n} \sum |X_i - \theta| \Rightarrow \theta_x$  median  $X$

To be always defined,  $\rho(x, \theta) := |x-\theta| - |x|$ . Then  $|E \rho(X, \theta)| = |E(|X-\theta| - |X|)| \leq E ||X-\theta| - |X|| \leq E|\theta| = \theta$  and  $E \rho(X, \theta)$  is defined even if  $E|X| = \infty$ .

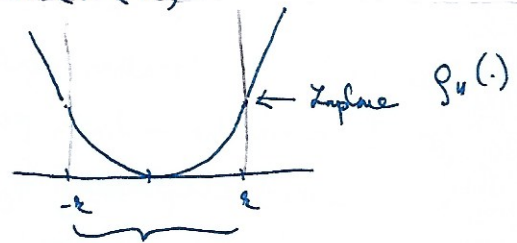
$\psi(x, \theta) = \frac{\partial}{\partial \theta} \rho(x, \theta) = \operatorname{sgn}(x-\theta)$  if  $x \neq \theta$  does not satisfy (25)

$\rho$  is convex in  $\theta \Rightarrow$  Theorem 10 gives  $\Gamma(\theta) = \frac{\partial^2}{\partial \theta^2} E \rho(X, \theta) \Big|_{\theta_x} = \dots = \frac{\partial^2 (2F(\theta) - 1)}{\partial \theta^2} \Big|_{\theta_x} = 2f(\theta_x)$

$\Sigma(\theta) = \operatorname{var} \psi(X, \theta_x) = 1$   $\Gamma_n(\hat{\theta} - \theta_x) \xrightarrow{d} N(0, \frac{1}{4f^2(F^{-1}(1/2))})$

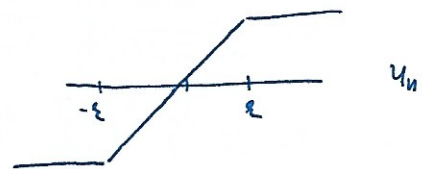
Huber's estimator  $\rho_h(x, \theta) = \begin{cases} \frac{(x-\theta)^2}{2} & |x-\theta| \leq h \\ h(|x-\theta| - \frac{h}{2}) & |x-\theta| > h \end{cases}$

$\psi_h(x, \theta) = \begin{cases} -(x-\theta) & |x-\theta| \leq h \\ -h \operatorname{sgn}(x-\theta) & |x-\theta| > h \end{cases}$



$\Sigma(\theta_x) = \operatorname{var} \psi_h(X, \theta_x) = E \psi_h(X, \theta_x)^2 = \int_{\theta_x-h}^{\theta_x+h} (x-\theta_x)^2 dF(x) + h^2 \int_{\theta_x+h}^{\infty} dF(x) + h^2 \int_{-\infty}^{\theta_x-h} dF(x)$

$\theta_x = \operatorname{argmin}_{\theta} E \rho_h(X, \theta)$  for  $f$  symmetric the median



$= \int_{\theta_x-h}^{\theta_x+h} (x-\theta_x)^2 dF(x) + h^2 (F(\theta_x-h) + 1 - F(\theta_x+h))$

$\Gamma(\theta) = \frac{\partial^2}{\partial \theta^2} E \rho_h(X, \theta) = \frac{\partial^2}{\partial \theta^2} \left( \int_{\theta-h}^{\theta+h} \frac{(x-\theta)^2}{2} dF(x) + h \int_{-\infty}^{\theta-h} (-(x-\theta) - \frac{h}{2}) dF(x) + h \int_{\theta+h}^{\infty} ((x-\theta) - \frac{h}{2}) dF(x) \right) =$

$= \frac{\partial}{\partial \theta} \left[ \frac{h^2}{2} f(\theta-h) - \frac{h^2}{2} f(\theta+h) + \int_{\theta-h}^{\theta+h} -(x-\theta) dF(x) + h \int_{-\infty}^{\theta-h} 1 dF(x) - h \int_{\theta+h}^{\infty} 1 dF(x) \right] =$

Leibniz integral rule

$\frac{\partial}{\partial \theta} \int_{a(\theta)}^{b(\theta)} f(x, \theta) dx = f(b(\theta), \theta) \frac{\partial b(\theta)}{\partial \theta} - f(a(\theta), \theta) \frac{\partial a(\theta)}{\partial \theta} + \int_{a(\theta)}^{b(\theta)} \frac{\partial f(x, \theta)}{\partial \theta} dx$   
 $dF(x) = f(x) dx$



$$= \frac{\partial}{\partial \theta} \left[ - \int_{\theta-h}^{\theta+h} (x-\theta) dF(x) + h F(\theta-h) - h (1-F(\theta+h)) \right] = -h f(\theta+h) - h f(\theta-h) + \int_{\theta-h}^{\theta+h} dF(x) + h f(\theta-h) + h f(\theta+h)$$

$$\frac{\partial}{\partial \theta} \int_{\theta-h}^{\theta+h} (x-\theta) dF(x) = h \cdot f(\theta+h) - (-h) f(\theta-h) + \int_{\theta-h}^{\theta+h} -1 dF(x) = F(\theta+h) - F(\theta-h)$$

**Theorem 10:**  $\Gamma_m(\hat{\theta} - \theta_x) \xrightarrow{d} N(0, \Gamma^{-1}(\theta_x) \Sigma(\theta_x) \Gamma^{-1}(\theta_x)) = N(0, \frac{\int_{\theta_x-h}^{\theta_x+h} (x-\theta_x)^2 dF(x) + h^2 (F(\theta_x+h) + 1 - F(\theta_x-h))}{(F(\theta_x+h) - F(\theta_x-h))^2})$

**M-estimation of location and scatter.**

$X_1, \dots, X_n \stackrel{iid}{\sim}$  location-scale model  $F_{\mu, \sigma} = F\left(\frac{x-\mu}{\sigma}\right)$   $\mu \in \mathbb{R}$  location parameter,  $\sigma > 0$  scale parameter, i.e. solve

MLE:  $f$  density of  $F$ , maximize  $\sum_{i=1}^n \log f\left(\frac{x_i-\mu}{\sigma}\right) - n \log \sigma$

$$\frac{1}{n} \sum_{i=1}^n \psi_1\left(\frac{x_i-\mu}{\sigma}\right) = 0 \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n \psi_2\left(\frac{x_i-\mu}{\sigma}\right) = 0 \quad \text{for} \quad \psi_1(t) = -\frac{f'(t)}{f(t)}, \quad \psi_2(t) = \psi_1(t)t - 1$$

Estimated parameters are identified by  $E \psi_i\left(\frac{X-\mu}{\sigma}\right) = 0 \quad i=1,2$ , influence function is proportional to  $\psi(x, \theta) = \begin{pmatrix} \psi_1\left(\frac{x-\mu}{\sigma}\right) \\ \psi_1\left(\frac{x-\mu}{\sigma}\right)\left(\frac{x-\mu}{\sigma}\right) - 1 \end{pmatrix}$ , asymptotic distribution follows from Theorems 9, 10.

Ex.  $N(\mu, \sigma^2)$   $f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$   $\psi_1(t) = t$   $\psi_2(t) = t^2 - 1$   $\hat{\mu} = \bar{X}$   $\hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{X})^2$

identify  $E \frac{X-\mu_x}{\sigma_x} = 0 \Rightarrow \mu_x = EX$ ,  $E \left(\frac{X-\mu_x}{\sigma_x}\right)^2 = 1 \Rightarrow \sigma_x^2 = \text{var} X$

influence function and asympt. normality derived already.

**Multivariate M-estimation of location and scatter.**

$X_1, \dots, X_n \stackrel{iid}{\sim}$  location-scale model with density  $f(x) = |\Sigma|^{-1/2} g\left(\left[(x-\mu)' \Sigma^{-1} (x-\mu)\right]^{1/2}\right)$  for  $\Sigma$  a scatter matrix, pos. def symmetric,  $\mu \in \mathbb{R}^d$  location parameter,  $g$  fixed function  $X_i$  is called elliptically symmetric with parameters  $g, \mu, \Sigma$ .

MLE:  $\ell_n(\mu, \Sigma) = -\frac{1}{2} n \log |\Sigma| + \sum_{i=1}^n \log g(\text{md}(x_i))$ ,  $\text{md}(x_i) = \left[(x_i-\mu)' \Sigma^{-1} (x_i-\mu)\right]^{1/2}$  Mahalanobis distance of  $x_i$  from  $\mu$ .

$V := \Sigma^{-1}$  and use Lehmann invariance principle

$$\frac{\partial \ell_n}{\partial \mu} = - \sum_{i=1}^n \frac{g'(\text{md}(x_i))}{g(\text{md}(x_i))} \frac{1}{2 \text{md}(x_i)} V(x_i-\mu) \stackrel{!}{=} 0 = \text{Tr}(\Sigma^{-1} (x_i-\mu)(x_i-\mu)^T)$$

$$\frac{\partial \ell_n}{\partial V} = -\frac{n}{2} (V^{-1})^T + \sum_{i=1}^n \frac{g'(\text{md}(x_i))}{g(\text{md}(x_i))} \frac{1}{2 \text{md}(x_i)} \frac{\partial}{\partial V} \text{Tr}[(x_i-\mu)' \Sigma^{-1} (x_i-\mu)]$$

$$= -\frac{n}{2} V^{-1} + \frac{1}{2} \sum_{i=1}^n \frac{g'(\text{md}(x_i))}{g(\text{md}(x_i))} \frac{1}{\text{md}(x_i)} (x_i-\mu)(x_i-\mu)^T \stackrel{!}{=} 0$$

We need  $\frac{\partial}{\partial V} \log |V| = (V^{-1})^T$   $\left[ V = \begin{pmatrix} a & c \\ c & d \end{pmatrix}, \log |V| = \log(ad-bc), \frac{\partial}{\partial V} \log |V| = \frac{1}{|V|} \begin{pmatrix} d & -c \\ -c & a \end{pmatrix} = \left[ \frac{1}{|V|} \begin{pmatrix} d & -c \\ -c & a \end{pmatrix} \right]^T = (V^{-1})^T \right]$

$\frac{\partial}{\partial V} \text{Tr}(V \cdot A) = \frac{\partial}{\partial V} \left( \sum_{i=1}^d (VA)_{ii} \right) = \frac{\partial}{\partial V} \sum_{i=1}^d \sum_{j=1}^d v_{ij} a_{ji} = A$  Set  $u(t) = -\frac{1}{t} \frac{\partial \log g(t)}{\partial t}$ . Then

we solve the system  $O_d = \frac{1}{n} \sum_{i=1}^n u(\text{md}(x_i)) (x_i-\mu)$  For  $u: [0, \infty) \rightarrow [0, \infty)$  arbitrary there define M-estimators of location and scatter.

$$O_{d \times d} = \frac{1}{n} \sum_{i=1}^n \left[ u(\text{md}(x_i)) (x_i-\mu)(x_i-\mu)^T - \Sigma \right]$$



**affine equivariance** an estimator of location is affine equivariant if for any  $A \in \mathbb{R}^{d \times d}$  non-singular and  $b \in \mathbb{R}^d$   $\mu_{AX+b} = A\mu_x + b$  for  $\mu_{AX+b}$  the estimator computed from data  $\{AX_i + b\}_{i=1}^m$ . An estimator of scatter  $\Sigma = \Sigma_x$  is affine equivariant if  $\Sigma_{AX+b} = A\Sigma_x A^T$ .

M-estimators of location and scatter are affine equivariant.

**Proof:**  $(AX+b - (A\mu_x + b))^T (A\Sigma_x A^T)^{-1} (AX+b - (A\mu_x + b)) = (x-\mu_x)^T A^T (A^T)^{-1} \Sigma_x^{-1} A (x-\mu_x) = (x-\mu_x)^T \Sigma_x^{-1} (x-\mu_x)$ . Thus  $E u(\text{md}(X)) (X-\mu_x) = 0$  if and only if  $E u(\text{md}(AX+b; A\mu_x + b, A\Sigma_x A^T)) (AX+b - (A\mu_x + b)) = A E u(\text{md}(X; \mu_x, \Sigma_x)) (X-\mu_x) = 0$ . Similarly  $E u(\text{md}(AX+b; A\mu_x + b, A\Sigma_x A^T)) (AX+b - (A\mu_x + b)) (AX+b - (A\mu_x + b))^T - A\Sigma_x A^T = A (E u(\text{md}(X; \mu_x, \Sigma_x)) (X-\mu_x)(X-\mu_x)^T - \Sigma_x) A^T = 0$  and  $\mu_x, \Sigma_x$  is an M-estimator for  $X \iff A\mu_x + b, A\Sigma_x A^T$  is an M-estimator for  $AX+b$   $\square$

In general, estimators of location without estimating scatter are not affine equivariant.

By affine equivariance it suffices to find M-estimators for  $\mu=0$  and  $\Sigma=I_d$ .

Influence functions of M-estimators of location and scatter are again proportional to  $\psi$ -function  $u(t) \cdot t$  and  $u(t) \cdot t^2$ . **Ex. 8.10.8.**

**(\*) Identification of parameters.** If the distribution of  $X_i$  is symmetric around 0 (WLOG  $\mu=0, \Sigma=I_d$ )

$E u(\text{md}(X)) (X-\mu_x) = 0$  is solved by  $\mu_x = 0$  [ $\text{md}(X) = \text{md}(-X)$ ]  $\Rightarrow \mu$  is correctly identified. If  $(X, Y) \stackrel{d}{=} (X, -Y)$  for  $X = (X, Y)^T$  and a bivariate case also the non-diagonal elements in  $E u(\text{md}(X)) X X^T - \Sigma_x$  are 0 and for the diagonal terms to make sure that  $E u(\text{md}(X)) X_i^2 - 1 = 0$  one has to consider  $\tilde{u}(t) = \frac{u(t)}{t}$  for  $u = E u(\text{md}(X)) X_i^2 = E u(\|X\|) X_i^2$ . Only then all parameters are correctly identified.

**Example: Cauchy distribution (or t-distribution with  $\nu$  degrees of freedom)**

$$g_\nu(t) = c \cdot (\nu + t^2)^{-\frac{\nu+1}{2}}$$

$$\log g_\nu(t) = \log c - \frac{\nu+1}{2} \log(\nu + t^2)$$

$$-\frac{\partial \log g_\nu(t)}{\partial t} = \frac{\nu+1}{2} \cdot \frac{2t}{\nu + t^2}$$

$$u(t) = \frac{d+\nu}{\nu+t^2}$$

influence functions are proportional to  $\frac{ct^j}{\nu+t^2}$ ,  $j=1,2$  and are all bounded.

asymptotic normality follows from Theorem 9 efficiency at normal model can be computed numerically. (Mathematica)

In practice, solutions to the estimating equations are obtained iteratively, or numerically.

**Identification of parameters for elliptically symmetric distributions**

$X \sim f(x) = |\Sigma|^{-1/2} g(\text{MD}(x))$ ,  $\text{MD}(x) := \sqrt{(x-\mu)^T \Sigma^{-1} (x-\mu)}$  we say  $X \sim \text{Ell}(\mu, \Sigma, g)$

Let  $b \in \mathbb{R}^d$  and  $A \in \mathbb{R}^{d \times d}$  positive definite. Then  $AX+b \sim \text{Ell}(A\mu+b, A\Sigma A^T, g)$ .

**Proof:** Transformation of a random vector:  $Y = AX+b$ ,  $X = A^{-1}(Y-b)$ ,  $|J| = |A^{-1}| = |A|^{-1}$



$$f_Y(y) = f_X(A'(y-b)) |A|^{-1} = |A|^{-1} |\Sigma|^{-1/2} g(\sqrt{(A'(y-b)-\mu)^T \Sigma^{-1} (A'(y-b)-\mu)}) =$$

$$= |A \Sigma A^T|^{-1/2} g(\sqrt{(y-(A\mu+b))^T (A')^T \Sigma^{-1} A' (y-(A\mu+b))}) \sim \text{Ell}(A\mu+b, A \Sigma A^T, g) \quad \square$$

In particular, if  $X \sim \text{Ell}(0, I, g) \Rightarrow AX \sim \text{Ell}(0, I, g)$  for any  $A$  orthogonal,  $X$  is invariant under orthogonal transformations (rotations)  $X$  is *spherically symmetric* i.e.  $AA^T = I$ .

**Identification:** By affine equivariance let  $\mu = 0, \Sigma = I, X \sim \text{Ell}(0, I, g)$ . Take  $A = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$

$$= \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \text{ A is OG and } AX = \begin{pmatrix} -x_1 \\ -x_2 \\ \vdots \\ -x_d \end{pmatrix} \stackrel{d}{=} \begin{pmatrix} x_1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} \Rightarrow \begin{matrix} X_1 \stackrel{d}{=} -X_1 \\ X \stackrel{d}{=} -X \end{matrix}$$

estimating equation of M-estimator gives  $E u(\|X\|) (X-0) = E(u(\|X\|)) (-X-0) = 0$

$$MD(X) = \sqrt{X^T I^{-1} X} = \|X\| \quad \text{and } \mu_X = 0 \text{ is correctly identified.}$$

$$\text{Further } E u(\|X\|) X X^T - I = \begin{cases} \text{if } j=k, j, k\text{-th element, } E u(\|X\|) X_j^2 - 1 = * \\ \text{if } j \neq k, E u(\|X\|) X_j X_k - 0 = E u(\|X\|) X_j (-X_k) = 0 \quad \checkmark \end{cases}$$

\* for  $*$   $\neq 0$  one has to assume that  $E u(\|X\|) X_j^2 = 1 \quad \forall j$ . Since  $X_j \stackrel{d}{=} X_1$  (take

$A = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$  I with  $j$ -th and 1-st  $\leftrightarrow$  now flipped) it suffices that  $E u(\|X\|) X_1^2 = 1 \Rightarrow$  modify and take the estimating equations with  $k = E u(\|X\|) X_j^2$

$$D_d = E u(MD(X)) (X - \mu)$$

$$D_{d \times d} = \frac{E u(MD(X)) (X - \mu) (X - \mu)^T - \Sigma}{E u(\|X\|) X_j^2} = E u_2(MD(X)) (X - \mu) (X - \mu)^T - \Sigma$$

for  $u_2(t) = u(t)/k$ .

Then both parameters are identified.

For  $d=1$   $0 = E u\left(\frac{|X-\mu|}{\sigma}\right) (X-\mu), \quad 0 = E \frac{u\left(\frac{|X-\mu|}{\sigma}\right) (X-\mu)^2}{E u(|X|) X^2} - \sigma^2$

$$MD(X) = \sqrt{\frac{(X-\mu)^2}{\sigma^2}} = \frac{|X-\mu|}{\sigma}$$

**M-estimators in regression.**

Ex. 38: Misspecified linear model.  $(\begin{smallmatrix} x_1 \\ y_1 \end{smallmatrix}), \dots, (\begin{smallmatrix} x_m \\ y_m \end{smallmatrix})$  iid,  $y_i | x_i \sim N(\beta' x_i, \sigma^2)$ , MLE  $\rightsquigarrow$  M-estimator

for  $\beta$  of form  $\rho(x_i, y_i; \beta) = (y - x' \beta)^2$

$$\psi(x_i, y_i; \beta) = -2x (y - x' \beta) \Rightarrow E X (Y - X' \beta_x) = 0 \Leftrightarrow E X Y = E [X X'] \beta_x$$

$$\beta_x = (E [X X'])^{-1} E X Y$$

if  $E [Y | X] = X' \beta \Rightarrow \beta_x = (E [X X'])^{-1} E X Y = (E [X X'])^{-1} E E [Y | X] = \beta$  and  $\beta$  is identified

$$E D_\psi(X, Y, \beta) = E 2 X X^T$$

$$E \psi(X, Y, \beta_x)^2 = 4 E X (Y - X' \beta_x) (X (Y - X' \beta_x))^T = 4 E E [X (Y - X' \beta_x) (Y - X' \beta_x)^T X^T | X] =$$

$$= 4 E X \text{diag}(var_{ix}(Y)) X^T = 4 E \sigma^2(X) X X^T$$

$\sigma^2(x_i) = var(Y | X = x_i)$

$$\text{var}(\hat{\beta} - \beta) \xrightarrow{d} N_p(0_p, (E X X^T)^{-1} (E \sigma^2(X) X X^T) (E X X^T)^{-1})$$

Sandwich estimator from linear regression.



**Robust regression estimators**

OLS:  $\rho(x, y, \beta) = (y - x'\beta)^2$

$\psi(x, y, \beta) = x(y - x'\beta)$

not robust, unbounded IF

LAD:  $\rho(x, y, \beta) = |y - x'\beta|$

$\psi(x, y, \beta) = x \operatorname{sign}(y - x'\beta)$

robust in  $Y$ , non-robust in  $X \rightarrow$  median regression

Huber:  $\rho(x, y, \beta) = \rho_H(y - x'\beta)$

$\psi(x, y, \beta) = \psi_H(y - x'\beta) \cdot x$

robust in  $Y$ , non-robust in  $X \rightarrow$  Huber's regression

**Identification of parameters. (OLS and LAD)**

Let  $Y = \alpha + X'\beta + \varepsilon$  for  $\varepsilon$  independent of  $X$ . Let the form of  $X$  be correct, i.e. the M-estimator

is given by  $\psi(x, y, a, b) = x(y - a - x'b)$ . Then the M-estimator identifies

OLS:

$E X(Y - a - X'b) = E X(\alpha + X'\beta + \varepsilon - a - X'b) = 0$

$\alpha EX + EXX'\beta + EXE\varepsilon - aEX - EXX'b = 0$

$EX(\alpha + E\varepsilon - a) + (EXX')(\beta - b) = 0$  and the choice  $a = \alpha + E\varepsilon$  solves the equation.  
 $b = \beta$

LAD: for given  $x$ ,  $\rho(x, y, a, b)$  is minimized by  $a + bx \approx$  conditional median of  $Y|X$ , which is  $\alpha + x'\beta + \operatorname{median}(\varepsilon)$ . This is a linear function of  $x \Rightarrow$  choice  $a = \alpha + \operatorname{med} \varepsilon$   
 $b = \beta$  identifies parameters.

~~for  $X$  given,  $E[Y - m|X]$  is minimized by  $m = \operatorname{median}(Y|X)$~~

In LAD, we fit a line  $m(x) = x'\beta$  that approximates  $\operatorname{median}(Y|X)$  the best in the sense  $E|Y - X'\beta|$  is minimized.

Example:  $Y = \alpha + \beta X + \varepsilon$ ,  $X \in \mathbb{R}$ ,  $X$  independent of  $\varepsilon$ . Misspecified model is  $Y \sim bX$ , no intercept in the model.

OLS:  $0 = E X(Y - Xb) = EX(\alpha + \beta X + \varepsilon - Xb) = \alpha EX + \beta EX^2 + EXE\varepsilon - bEX^2 = 0$

$b = \frac{1}{EX^2} (\alpha EX + \beta EX^2 + EXE\varepsilon) = \beta + \frac{\alpha + E\varepsilon}{EX^2} \cdot EX$  [ in R script  $\alpha = 0, EX = 1/2, EX^2 = 1/3$   
 $E\varepsilon = 1$  (Exp(1)) and  $b = 5/2$   
 $\beta = 1$  ]

conditional expectation:  $E[Y|X] = \alpha + E\varepsilon + \beta X$

~~$E[(Y - \frac{Y + m(x)}{2})^2 | X] = E[(Y - E[Y|X])^2 | X] + E[(E[Y|X] - m(x))^2 | X]$~~

~~$+ 2E[(Y - E[Y|X])(E[Y|X] - m(x)) | X]$~~

~~$= 0$~~

and for given  $X$ ,  $E \rho(x, Y, b)$  is minimized iff  $m(x)$  is close in the  $L^2$ -norm ( $X$ )

to  $E(Y|X)$   ~~$(Y - m(x))^2 \Rightarrow$  it is enough to minimize  $E((E[Y|X] - m(x))^2)$~~

Pythagorean theorem: Hilbert space  $L^2(X)$  given by  $\|f\|_X^2 = E f(X)^2$

LAD: minimize  $E| \alpha + \beta X + \varepsilon - bX | = E E[| \alpha + \beta X + \varepsilon - bX | | X]$

for any  $X$  the inner expectation is minimized by conditional median  $\alpha + \beta X + \operatorname{med} \varepsilon$  but this is not a function of the form  $bX$  and because  $L^1(X)$  is not Hilbert, no version Pythagorean theorem can be stated. We do not minimize  $E|m(x) - \operatorname{med}(Y|X)|$

OLS:  $\rho(x, y) := (y - m(x))^2$  minimize over functions  $m \in \mathcal{M}$

$$E \rho(x, y) = E (y - m(x))^2 = E E[(y - m(x))^2 | x] = E \left( E[(y - E(y|x))^2 | x] + E[(E(y|x) - m(x))^2 | x] \right) \\ + 2 E \left[ (y - E(y|x))(E(y|x) - m(x)) | x \right] = \underbrace{E (y - E(y|x))^2}_{\text{does not depend}} + \underbrace{E (E(y|x) - m(x))^2}_{= \|E(y|x) - m\|_x^2}$$

In OLS, always the closest function to  $E(y|x)$  in the  $\|\cdot\|_x$ -norm is specified.

This holds because of the Hilbert-space structure of  $\|\cdot\|_{x,2}$  in  $L^2(X)$ ,  $\|f\|_{x,2} = \sqrt{E f(x)^2}$

No such thing holds true in the for  $L^1(X)$ , LAD does not specify the closest function to median  $(y|x)$  in the  $L^1(X)$ -norm.  $\|f\|_{x,1} = E |f(x)|$



Example 3.9:  $x_i = (x_{i1}, \dots, x_{ip})^T$   $y_i | x_i \sim P_o(\lambda(x_i))$

$$\lambda(x_i) = e^{\beta x_i} \quad \beta = (\beta_1, \dots, \beta_p)^T$$

$$L(\beta) = \prod \frac{e^{-\lambda(x_i)} \lambda(x_i)^{y_i}}{y_i!} \quad \ell(\beta) = -\sum \lambda(x_i) + \sum y_i \log \lambda(x_i) + c$$

$$\frac{\partial}{\partial \beta} \lambda(x_i) = x_i \lambda(x_i) \quad U(\beta) = -\sum x_i \lambda(x_i) + \sum y_i x_i = \sum x_i (y_i - \lambda(x_i))$$

$\Rightarrow \psi(x_i, y_i, \beta) = x_i (y_i - e^{\beta x_i})$  and  $\beta_x$  solves (the true parameter)

$$E \psi(x_i, y_i, \beta_x) = 0$$

Interoctually, if different distribution

If  $y_i | x_i$  is not Poisson, but still  $E[y_i | x_i] = \lambda(x_i)$  we have for  $\beta_0 = \beta$

$$E \psi(x_i, y_i, \beta_0) = E X (Y - \lambda(X)) = E X (Y - e^{\beta_0 X})$$

$$= E E [X (Y - e^{\beta_0 X}) | X] = 0 \quad \Rightarrow \quad \underline{\underline{\beta_0 = \beta_x}}$$

and we still estimate the correct parameter.

Theorem 9:  $\Pi(\beta_x) = E \frac{\partial}{\partial \beta} \psi(x_i, y_i, \beta) |_{\beta_x} = E X X' e^{\beta_x X}$

$$\Sigma(\beta_x) = E X (Y - e^{\beta_x X})^2 X' = E X X' (Y - e^{\beta_x X})$$

$$\Gamma_n(\hat{\beta} - \beta_x) \xrightarrow{d} N_p(0, \Pi^{-1}(\beta_x) \Sigma(\beta_x) \Pi^{-1}(\beta_x))$$

estimate  $\Pi$  a  $\Sigma$  as usual from the observed quantities

$\Rightarrow$  Sandwich estimator in Poisson regression

If the true Poisson is true, even is given by Fisher information

$\sim \Sigma$  and  $\Pi$  (both correspond to  $I(\beta)$ )