Lecture Stat 461-561 M-Estimation

AD

February 2008

- In most applications, we have X_i ~ g and we obtain an estimate a by minimizing a suitable cost function; e.g.
 - the mean corresponds to $\sum_{i=1}^{n} \left(\theta x_i \right)^2$.
 - the median corresponds to $\sum_{i=1}^{n} |\theta x_i|$.
 - the MLE corresponds to negative log-likelihood $-\sum_{i=1}^{n} \log f(x_i | \theta)$.
- However, even the mean estimate of a location parameter is typically not robust. In contrast, the median could be too 'rough'.
- Example: Consider

 $\mathbf{x} = (-1.28, -0.96, -0.46, -0.44, -0.26, -0.21, -0.063, 0.39, 3, 6, 9)$

where $X_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ for i = 1, ..., 8 but X_9, X_{10}, X_{11} are outliers.

- The mean is 1.33 and the median is -0.21.
- Huber (1964) introduced a general loss function which is a compromise between mean and median; i.e. we minimize

$$\sum_{i=1}^{n} \rho\left(x_i - \theta\right)$$

where

$$\rho(x) = \begin{cases} \frac{1}{2}x^2 & \text{if } |x| \le k \\ k |x| - \frac{1}{2}k^2 & \text{if } |x| \ge k \end{cases}$$

- Large observations are not as heavily weighted as for $\sum_{i=1}^{n} (\theta x_i)^2$.
- *k* is a tuning parameter which controls the mix between the mean and median-like estimators.

Results

k	0	1	2	3	4	5	6	8	10
Estimate	21	.03	04	.29	.41	.52	.87	.97	1.33

Huber's estimator as a function of k

- When k = 0, Huber's estimator corresponds to the median and as k increases it gets closer to the mean; i.e. the robustness properties of the estimator are decreasing.
- Remark: Clearly minimizing

$$\sum_{i=1}^{n} \rho\left(x_i - \theta\right)$$

appears equivalent to maximizing a log-likelihood for which $\log f(x_i | \theta) = cste - \rho(x_i - \theta)$. We will describe later more general estimates which do not support such an interpretation.

• Now we assume a more general case where $X_i \sim g$ and our estimate $\hat{\theta}_n$ is solution of

$$\sum_{i=1}^{n}\psi(x_{i},\theta)=0.$$

• Under regularity conditions $\widehat{\theta}_n$ will converge towards the parameter θ^* satisfying

$$\mathbb{E}_{g}\left[\psi\left(X,\theta^{*}\right)\right]=\int\psi\left(x,\theta^{*}\right)g\left(x\right)dx=0.$$

• If $g(x) = f(x|\theta_0)$ then θ^* is defined by

$$\mathbb{E}_{f\left(\cdot\mid\theta\right)}\left[\psi\left(X,\theta^{*}\right)\right]=\int\psi\left(x,\theta^{*}\right)f\left(x\mid\theta_{0}\right)dx=0,$$

i.e. be careful: we do not have necessarily $\theta^* = \theta_0!$ Also in practice, we would like it to be the case.

• For example, if $\psi(x, \theta) = x - \theta$ then $\theta^* = \mathbb{E}_g[X]$.

• To study the asymptotic properties of $\hat{\theta}_n$, we use the (now standard) Taylor expansion method of $\sum_{i=1}^n \psi(\theta, x_i)$ around the value θ^* which yields

$$0 = \sum_{i=1}^{n} \psi(x_{i}, \theta^{*}) + \left(\widehat{\theta}_{n} - \theta^{*}\right) \sum_{i=1}^{n} \psi'(x_{i}, \theta^{*}) + R_{n}$$

• By ignoring the remainder term R_n , we obtain

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta^{*}\right)=\frac{-\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\psi\left(x_{i},\theta^{*}\right)}{\frac{1}{n}\sum_{i=1}^{n}\psi'\left(x_{i},\theta^{*}\right)}$$

• The CLT yields

$$-\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\psi\left(x_{i},\theta^{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(0,\mathbb{E}_{g}\left[\psi^{2}\left(X,\theta^{*}\right)\right]\right)$$

as $\mathbb{E}_{g} \left[\psi \left(X_{i}, \theta^{*} \right) \right] = 0$ and $var_{g} \left[\psi \left(X_{i}, \theta^{*} \right) \right] = \mathbb{E}_{g} \left[\psi^{2} \left(X, \theta^{*} \right) \right]$. • The law of large numbers provides

$$\frac{1}{n}\sum_{i=1}^{n}\psi'\left(x_{i},\theta^{*}\right)\xrightarrow{\mathsf{P}}\mathbb{E}_{g}\left[\psi'\left(X,\theta^{*}\right)\right]$$

So by Slutsky's theorem, we get

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta^{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(\mathsf{0},\frac{\mathbb{E}_{g}\left[\psi^{2}\left(X,\theta^{*}\right)\right]}{\mathbb{E}_{g}\left[\psi^{\prime}\left(X,\theta^{*}\right)\right]^{2}}\right)$$

 This is a generalization of the misspecified model we discussed before where ψ is arbitrary.

Study of the Huber's estimate

• For the Huber's estimate, we have $\psi\left(x, heta
ight)=\psi\left(x- heta
ight)$ where

$$\psi(x) = \begin{cases} x & \text{if } |x| \le k \\ k & \text{if } x \ge k \\ -k & \text{if } x < -k \end{cases}$$

 Assume we have X_i ^{i.i.d.} r (x − θ) where f is symmetric around 0 and we want to estimate θ then indeed

$$\mathbb{E} \left[\psi \left(X - \theta \right) \right] \\= \int_{\theta - k}^{\theta + k} \left(x - \theta \right) f \left(x - \theta \right) dx - k \int_{-\infty}^{\theta - k} f \left(x - \theta \right) dx \\+ k \int_{\theta + k}^{+\infty} f \left(x - \theta \right) dx \\= \int_{-\infty}^{k} uf \left(u \right) du - k \int_{-\infty}^{-k} f \left(u \right) du + k \int_{k}^{+\infty} f \left(u \right) du \\= 0.$$

• In this case we have indeed that $\theta^* = \theta!$

AD ()

• We also have

$$\mathbb{E}\left[\psi'\left(X- heta
ight)
ight]=\int_{ heta-k}^{ heta+k}f\left(x- heta
ight)dx=P_{ heta}\left(|X|\leq k
ight),$$

$$\mathbb{E}\left[\psi^{2}\left(X-\theta\right)\right] = \int_{\theta-k}^{\theta+k} \left(x-\theta\right)^{2} f\left(x-\theta\right) dx + k^{2} \int_{-\infty}^{\theta-k} f\left(x-\theta\right) dx \\ + k^{2} \int_{\theta+k}^{+\infty} f\left(x-\theta\right) dx \\ = \int_{-k}^{k} u^{2} f\left(u\right) du + 2k^{2} \int_{k}^{+\infty} f\left(u\right) du.$$

• It follows that the Huber's estimate satisfies

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta^{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(0,\frac{\int_{-k}^{k}u^{2}f\left(u\right)du+2k^{2}P_{\theta^{*}}\left(|X|>k\right)}{\left[P_{\theta^{*}}\left(|X|\leq k\right)\right]^{2}}\right)$$

• We compare the asymptotic relative efficiencies of Huber's estimate for k = 1.5 to mean and median

	Normal	Double Exponential		
vs. mean	.96	1.37		
vs. median	1.51	.68		

that is
$$\sigma^2_{
m Huber}/\sigma^2_{
m mean}$$
 and $\sigma^2_{
m Huber}/\sigma^2_{
m median}.$

- Remember that mean is the MLE of normal and median is the MLE of double exponential so ARE are <1 as expected.
- Huber's estimator performs however reasonably well compared to the MLE.

- An M-estimator is a tradeoff between robustness and efficiency.
- To see how much we are losing, we study in more details the asymptotic variance given by $\mathbb{E}\left[\psi'\left(X,\theta^*\right)\right]^{-2}\mathbb{E}\left[\psi^2\left(X,\theta^*\right)\right]$.
- We have

$$\mathbb{E}\left[\psi'\left(X,\theta\right)\right] = -\int \frac{d\psi\left(x,\theta\right)}{d\theta}f\left(x|\theta\right)dx$$

where

$$\frac{d}{d\theta} \int \psi(x,\theta) f(x|\theta) dx = \int \frac{d\psi(x,\theta)}{d\theta} f(x|\theta) dx + \int \psi(x,\theta) \frac{df(x|\theta)}{d\theta} dx$$

so if $\int \psi(x,\theta) f(x|\theta) dx = 0$ for all θ then

$$\mathbb{E}\left[\psi'\left(X,\theta\right)\right] = \int \psi\left(x,\theta\right) \frac{df\left(x|\theta\right)}{d\theta} dx$$
$$= \int \psi\left(x,\theta\right) \frac{d\log f\left(x|\theta\right)}{d\theta} f\left(x|\theta\right) dx$$

Recall that the asymptotic variance of the MLE is in

$$\mathbb{E}_{\theta}\left[\frac{d\log f(X|\theta)}{d\theta}^2\right]^{-1}$$
 thus

$$ARE = \frac{var\left(\mathsf{MLE}\right)}{var\left(\mathsf{M}\right)} = \frac{\mathbb{E}\left[\psi\left(X,\theta\right)\frac{d\log f\left(X|\theta\right)}{d\theta}\right]^{2}}{\mathbb{E}\left[\psi^{2}\left(X,\theta^{*}\right)\right]\mathbb{E}_{\theta}\left[\frac{d\log f\left(X|\theta\right)}{d\theta}^{2}\right]} \le 1$$

follows from the Cauchy-Schwartz inequality.

- An M-estimate is always less efficient than the MLE and matches its efficiency only if $\psi(x, \theta)$ is proportional to $\frac{d \log f(x|\theta)}{d\theta}$.
- This result does not say much; if one uses an M-estimate it is because it it not believed that the model $f(x|\theta)$ is reliable...

General multivariate case

• We want to estimate the multidimensional parameter θ^* which satisfies

$$\mathbb{E}\left[\psi\left(X,\theta^*\right)\right] = \int \psi\left(x,\theta^*\right) f\left(x\right) dx = 0.$$

• This extension is trivial theoretically but will allow us to study numerous interesting estimates; e.g. consider the estimate

$$\widehat{\theta}_{1,n} = \frac{1}{n} \sum_{i=1}^{n} |x_i - \overline{x}|.$$

• At first glance, this is not an M-estimate as there is no single equation of the form

$$\sum_{i=1}^{n}\psi(x_{i},\theta_{1})=0$$

that yields $\widehat{\theta}_{1,n}$.

Moreover there is no family of densities f (x| θ) such that θ
_{1,n} is a component of the MLE of θ.

• However, we can write

$$\psi(\mathbf{x}, \theta) = \begin{pmatrix} \psi_1(\mathbf{x}, \theta_1, \theta_2) \\ \psi_2(\mathbf{x}, \theta_1, \theta_2) \end{pmatrix} = \begin{pmatrix} |\mathbf{x} - \theta_2| - \theta_1 \\ (\mathbf{x} - \theta_2) \end{pmatrix}$$

• We find out that

$$\sum_{i=1}^{n}\psi(x_i,\theta_1,\theta_2)=0$$

implies $\hat{\theta}_{2,n} = \frac{1}{n} \sum_{i=1}^{n} x_i$, $\hat{\theta}_{1,n} = \frac{1}{n} \sum_{i=1}^{n} |x_i - \overline{x}|$.

n

• Asymptotic results can be easily establish using a straightforward generalization of the scalar case and, under regularity assumptions, we obtain

$$\sqrt{n}\left(\widehat{\theta}_{n}-\theta^{*}\right)\xrightarrow{\mathsf{D}}\mathcal{N}\left(0,V\left(\theta^{*}\right)\right)$$

where

$$V\left(\theta^{*}\right) = A^{-1}\left(\theta^{*}\right) B\left(\theta^{*}\right) \left\{A^{-1}\left(\theta^{*}\right)\right\}^{\mathsf{I}}$$

with

$$A(\theta^*) = \mathbb{E}\left[-\frac{\partial \psi(X,\theta)}{\partial \theta^{\mathsf{T}}}\right]\Big|_{\theta=\theta^*}, \ B(\theta^*) = \mathbb{E}\left[\psi(X,\theta^*)\psi(X,\theta^*)^{\mathsf{T}}\right]$$

• Clearly if $\psi(x, \theta) = \frac{\partial \log f(x|\theta)}{\partial \theta}$ and if the data truly come from the assumed parametric family $f(x|\theta)$ then

$$A(\theta^*) = B(\theta^*) = I(\theta^*) \Rightarrow V(\theta^*) = I(\theta^*)^{-1}.$$

• However, in many cases the data do not come from the assumed family and valid inference should be carried out using the correct limiting covariance matrix.

• Let us define $G_n(\theta) := \sum_{i=1}^n \psi(x_i, \theta)$. The idea of the proof is always the same

$$G_{n}\left(\widehat{\theta}_{n}\right)=0=G_{n}\left(\theta^{*}\right)+G_{n}^{\prime}\left(\theta^{*}\right)\left(\widehat{\theta}_{n}-\theta^{*}\right)+R_{n}$$

where
$$G'_{n}(\theta^{*}) = \left[\frac{\partial G_{n}(\theta)}{\partial \theta^{T}}\right]\Big|_{\theta=\theta^{*}}$$
 so
 $\sqrt{n}\left(\widehat{\theta}_{n}-\theta^{*}\right) = -\left[G'_{n}(\theta^{*})\right]^{-1}\sqrt{n}G_{n}(\theta^{*}) + \sqrt{n}R_{n}^{*}$

• Moreover under regularity conditions

$$\begin{split} & -G'_{n}\left(\theta^{*}\right) \xrightarrow{\mathsf{P}} A\left(\theta^{*}\right), \\ & \sqrt{n}G_{n}\left(\theta^{*}\right) \xrightarrow{\mathsf{D}} \mathcal{N}\left(0, B\left(\theta^{*}\right)\right), \\ & \sqrt{n}R_{n}^{*} \xrightarrow{\mathsf{P}} 0. \end{split}$$

• We can estimate $A\left(\theta^{*}\right)$, $B\left(\theta^{*}\right)$ and $V\left(\theta^{*}\right)$ using the data samples via

$$\begin{aligned} A_n\left(\widehat{\theta}_n, \mathbf{x}\right) &= \frac{1}{n} \sum_{i=1}^n -\frac{\partial \psi\left(x_i, \widehat{\theta}_n\right)}{\partial \theta^{\mathsf{T}}}, \\ B\left(\widehat{\theta}_n, \mathbf{x}\right) &= \frac{1}{n} \sum_{i=1}^n \psi\left(x_i, \widehat{\theta}_n\right) \psi\left(x_i, \widehat{\theta}_n\right)^{\mathsf{T}}, \\ V\left(\widehat{\theta}_n, \mathbf{x}\right) &= A^{-1}\left(\widehat{\theta}_n, \mathbf{x}\right) B\left(\widehat{\theta}_n, \mathbf{x}\right) \left\{A^{-1}\left(\widehat{\theta}_n, \mathbf{x}\right)\right\}^{\mathsf{T}} \end{aligned}$$

• Under mild regularity assumptions, we have

$$V\left(\widehat{\theta}_{n},\mathbf{x}\right)\xrightarrow{\mathsf{P}}V\left(\theta^{*}\right)$$

An interesting extension consists of considering

$$\sum_{i=1}^{n}\psi\left(x_{i}, heta
ight)=c_{n}$$

where $c_n / \sqrt{n} \xrightarrow{\mathsf{P}} 0$ as $n \to \infty$.

• In these cases, the asymptotic results still hold as we can simply write

$$G_{n}\left(\widehat{\theta}_{n}\right) - c_{n} = 0 = G_{n}\left(\theta^{*}\right) - c_{n} + G_{n}'\left(\theta^{*}\right)\left(\widehat{\theta}_{n} - \theta^{*}\right) + R_{n}$$
$$= G_{n}\left(\theta^{*}\right) + G_{n}'\left(\theta^{*}\right)\left(\widehat{\theta}_{n} - \theta^{*}\right) + R_{n} - c_{n}$$

and c_n is absorbed in the remainder.

• *Example*. **Posterior mode**. In this case, assume we are interested in maximizing the posterior distribution which is proportional to

$$\pi\left(\theta\right)\prod_{i=1}^{n}f\left(\left.x_{i}\right|\theta\right)$$

Then it can be written as

$$\sum_{i=1}^{n} \frac{\partial \log f\left(\left.x_{i}\right|\theta\right)}{\partial \theta} = -\frac{\partial \log \pi\left(\theta\right)}{\partial \theta}.$$

• It follows that as long as

$$c_{n}\left(heta
ight)=-rac{\partial\log\pi\left(heta
ight)}{\partial heta}$$

is such that $c_n(\theta) / \sqrt{n} \xrightarrow{P} 0$ then the Bayesian MAP estimator has the same asymptotic covariance as the MLE.

• Let $\hat{\theta}_n = (\overline{x}_n, s_n) = \left(\frac{1}{n} \sum_{i=1}^n x_i, \frac{1}{n} \sum_{i=1}^n (x_i - \overline{x}_n)^2\right)$. This estimate is an M-estimate for

$$\psi(x,\theta) = \begin{pmatrix} \psi_1(x,\theta_1,\theta_2) \\ \psi_2(x,\theta_1,\theta_2) \end{pmatrix} = \begin{pmatrix} x-\theta_1 \\ (x-\theta_1)^2 - \theta_2 \end{pmatrix}.$$

We can calculate

$$\begin{split} A\left(\theta^{*}\right) &= \mathbb{E}\left[-\frac{\partial\psi(X,\theta)}{\partial\theta^{\mathsf{T}}}\right]\Big|_{\theta=\theta^{*}} \\ &= \mathbb{E}\left(\begin{array}{cc}1&0\\2\left(x-\theta_{1}^{*}\right)&1\end{array}\right) = \left(\begin{array}{cc}1&0\\0&1\end{array}\right), \\ B\left(\theta^{*}\right) &= \mathbb{E}\left[\psi\left(X,\theta^{*}\right)\psi\left(X,\theta^{*}\right)^{\mathsf{T}}\right] \\ &= \mathbb{E}\left(\begin{array}{cc}\left(x-\theta_{1}^{*}\right)^{2}&\left(x-\theta_{1}^{*}\right)\left(\left(x-\theta_{1}^{*}\right)^{2}-\theta_{2}^{*}\right)\\\left(x-\theta_{1}^{*}\right)\left(\left(x-\theta_{1}^{*}\right)^{2}-\theta_{2}^{*}\right)&\left(\left(x-\theta_{1}^{*}\right)^{2}-\theta_{2}^{*}\right)^{2}\\ &= \left(\begin{array}{cc}\theta_{2}^{*}&\mu_{3}\\\mu_{3}&\mu_{4}-\theta_{2}^{*2}\end{array}\right) = \left(\begin{array}{cc}\theta_{2}^{*}&\mu_{3}\\\mu_{3}&\mu_{4}-\sigma^{4}\end{array}\right) \end{split}$$

where μ_3 is the 3th central moment of X and we have use the more familiar notation $\sigma^2 = \theta_2^*$.

• We can estimate $B\left(heta ^{st }
ight)$ by

$$B\left(\widehat{\theta}_n,\mathbf{x}\right) = \begin{pmatrix} s_n^2 & m_3 \\ m_3 & m_4 - s_n^4 \end{pmatrix}.$$

•
$$\hat{\theta}_n$$
 is a MLE estimate associated to
 $f(x|\theta) = (2\pi\theta_2)^{-1/2} \exp\left(-(x-\theta_1)^2/2\theta_2\right)$ but
 $\psi_1(x,\theta_1,\theta_2) = x-\theta_1$ and $\psi_2(x,\theta_1,\theta_2) = (x-\theta_1)^2-\theta_2$ are not
the score functions which are equal to $\frac{\partial \log f(x|\theta)}{\partial \theta_1} = (x-\theta_1)/\theta_2$ and
 $\frac{\partial \log f(x|\theta)}{\partial \theta_2} = (x-\theta_1)^2/2\theta_2^2 - 1/2\theta_2.$

• It follows that clearly the ψ functions are not unique - many different functions lead to the same estimator. They also yield different $A(\theta^*)$ and $B(\theta^*)$ but the same $V(\theta^*)$.

• If we pick $\psi_{\mathsf{MLE}}\left(\mathbf{x},\theta\right)=\frac{\partial\log f\left(\mathbf{x}|\theta\right)}{\partial\theta}$ then

$$\begin{aligned} A\left(\theta^{*}\right) &= \mathbb{E}\left[-\frac{\partial\psi\left(X,\theta\right)}{\partial\theta^{\mathsf{T}}}\right]\Big|_{\theta=\theta^{*}} = \begin{pmatrix} 1/\sigma^{2} & 0\\ 0 & 1/\sigma^{4} \end{pmatrix}, \\ B\left(\theta^{*}\right) &= \mathbb{E}\left[\psi\left(X,\theta^{*}\right)\psi\left(X,\theta^{*}\right)^{\mathsf{T}}\right] = \begin{pmatrix} 1/\sigma^{2} & \frac{\mu_{3}}{2\sigma^{3}}\\ \frac{\mu_{3}}{2\sigma^{3}} & \frac{\mu_{4}-\sigma^{4}}{4\sigma^{8}} \end{pmatrix} \end{aligned}$$

• If the data are distributed according to $f(x|\theta)$ then $\mu_3 = 0$ and $\mu_4 = 3\sigma^4$ and it follows that

$$A\left(heta^{*}
ight)=B\left(heta^{*}
ight)= extsf{Diag}\left(1/\sigma^{2},1/\sigma^{4}
ight).$$

• Note that the likelihood score functions $\psi_{\rm MLE}$ are related to the original ψ by

$$\psi_{\mathsf{MLE}} = \mathcal{C}\psi$$

where $C = Diag (1/\sigma^2, 1/\sigma^4)$. Generally speaking all functions $\psi' = C\psi$ where C is non singular (but possibly dependent on θ^* and **x**) leads to the same estimator and the same asymptotic matrix.

• Example Ratio Estimator: Let

$$\widehat{\theta}_n = \frac{\overline{y}}{\overline{x}}$$

where
$$\overline{x} = n^{-1} \sum_{i=1}^{n} x_i$$
 and $\overline{y} = n^{-1} \sum_{i=1}^{n} y_i$ with $\mathbb{E}(X) = \mu_X$, $\mathbb{E}(Y) = \mu_Y$, $var(X) = \sigma_X^2$, $var(Y) = \sigma_Y^2$ and $cov(X, Y) = \sigma_{XY}$.
• We have

$$\psi(X, Y, \theta) = Y - \theta X$$

thus

$$\begin{array}{lll} A\left(\theta^{*}\right) & = & \mathbb{E}\left[-\frac{\partial\psi\left(X,\theta\right)}{\partial\theta^{\mathsf{T}}}\right]\Big|_{\theta=\theta^{*}} = \mu_{X}, \\ B\left(\theta^{*}\right) & = & \mathbb{E}\left[\psi\left(X,\theta^{*}\right)^{2}\right] = \mathbb{E}\left[\left(Y-\theta X\right)^{2}\right], \\ V\left(\theta^{*}\right) & = & \mathbb{E}\left[\left(Y-\theta^{*}X\right)^{2}\right]/\mu_{X}^{2}, \end{array}$$

• These matrices can be estimated through

$$\begin{array}{lll} A\left(\widehat{\theta}_{n},\mathbf{x},\mathbf{y}\right) &=& \overline{x}, \\ B\left(\widehat{\theta}_{n},\mathbf{x},\mathbf{y}\right) &=& \displaystyle\frac{1}{n}\sum_{i=1}^{n}\left(y_{i}-\frac{\overline{y}}{\overline{x}}x_{i}\right)^{2}, \\ V\left(\widehat{\theta}_{n},\mathbf{x},\mathbf{y}\right) &=& \displaystyle\frac{1}{n\overline{x}^{2}}\sum_{i=1}^{n}\left(y_{i}-\frac{\overline{y}}{\overline{x}}x_{i}\right)^{2}. \end{array}$$

• If we are interested in the joint distribution of $(\overline{x}, \overline{y}, \frac{\overline{y}}{\overline{x}})$, we only need to define

$$\psi(X, Y, \theta) = \begin{pmatrix} Y - \theta_1 \\ X - \theta_2 \\ \theta_1 - \theta_3 \theta_2 \end{pmatrix}$$

We obtain

$$A(\theta^*) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & \theta_3^* & \theta_2^* \end{pmatrix}, \ B(\theta^*) = \begin{pmatrix} \sigma_Y^2 & \sigma_{XY} & 0 \\ \sigma_{XY} & \sigma_X^2 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

• We can check that the (3,3)th element of $V(\theta^*) = A^{-1}(\theta^*) B(\theta^*) \{A^{-1}(\theta^*)\}^{\mathsf{T}}$ is

$$v_{33} = \frac{1}{\theta_2^{*2}} \left[\sigma_Y^2 - 2\theta_3^* \sigma_{XY} + \theta_3^{*2} \sigma_X^2 \right]$$
$$= \mathbb{E} \left[\left(\mathbf{Y} - \theta^* \mathbf{X} \right)^2 \right] / \mu_X^2$$

• Example Instrumental Variable Estimation:

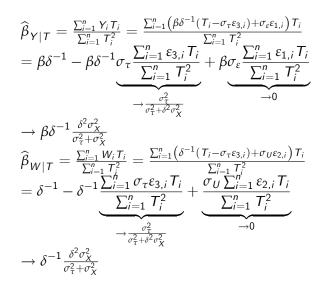
$$\begin{array}{lll} Y_i &=& \beta X_i + \sigma_{\varepsilon} \varepsilon_{1,i} \\ W_i &=& X_i + \sigma_U \varepsilon_{2,i} \\ T_i &=& \gamma + \delta X_i + \sigma_{\tau} \varepsilon_{3,i} \end{array}$$

where $\varepsilon_{j,i}$ are mutually independent erros with zero mean and unit variance. We also assume that $X_1, ..., X_n$ are unobserved, independent of $\{\varepsilon_{j,i}\}$ and have finite variance σ_X^2 .

- W_i is a measurement of X_i and T_i is an instrumental variable for X_i (for estimating β) provided that δ ≠ 0.
- The OLS estimator of slope obtained by regressing Y on W is

$$\begin{split} \widehat{\boldsymbol{\beta}}_{Y|W} &= \frac{\sum_{i=1}^{n} W_{i}Y_{i}}{\sum_{i=1}^{n} W_{i}^{2}} = \frac{\sum_{i=1}^{n} W_{i}(\boldsymbol{\beta}(W_{i} - \sigma_{U}\boldsymbol{\epsilon}_{2,i}) + \sigma_{\boldsymbol{\epsilon}}\boldsymbol{\epsilon}_{1,i})}{\sum_{i=1}^{n} W_{i}^{2}} \\ &= \boldsymbol{\beta} - \boldsymbol{\beta}\underbrace{\frac{\sigma_{U}\sum_{i=1}^{n} W_{i}}{\sum_{i=1}^{n} W_{i}^{2}}}_{\rightarrow \frac{\sigma_{U}^{2}}{\sigma_{U}^{2} + \sigma_{X}^{2}}} + \underbrace{\frac{\sigma_{\boldsymbol{\epsilon}}\sum_{i=1}^{n} W_{i}\boldsymbol{\epsilon}_{1,i}}{\sum_{i=1}^{n} W_{i}^{2}}}_{\rightarrow 0} \\ &\stackrel{\mathbf{P}}{\rightarrow} \frac{\sigma_{X}^{2}}{\sigma_{X}^{2} + \sigma_{U}^{2}}\boldsymbol{\beta}. \end{split}$$

• For sake of simplicity, lets take here $\gamma = 0$. Let $\hat{\beta}_{Y|W}$ and $\hat{\beta}_{W|T}$ be the slopes from the LS regressions of Y on T and W. We have



• The intrumental variable estimator is defined by

$$\widehat{\beta}_{\mathsf{IV}} = \frac{\widehat{\beta}_{Y|T}}{\widehat{\beta}_{W|T}} = \frac{\sum_{i=1}^{n} Y_i T_i}{\sum_{i=1}^{n} W_i T_i} \to \beta.$$

ullet This estimate is an M-estimator. A choice for ψ consists of using

$$\psi(Y, W, T, \theta) = \begin{pmatrix} \theta_1 - T \\ (Y - \theta_2 W)(\theta_1 - T) \end{pmatrix}.$$

Indeed

$$\frac{1}{n}\sum_{i=1}^{n} (\theta_1 - t_i) = 0 \Rightarrow \widehat{\theta}_1 = \overline{t} = \frac{1}{n}\sum_{i=1}^{n} t_i,$$
$$\frac{1}{n}\sum_{i=1}^{n} (y_i - \theta_2 w_i) (\theta_1 - t_i) = 0 \Rightarrow \widehat{\theta}_2 = \frac{\sum_{i=1}^{n} y_i (t_i - \overline{t})}{\sum_{i=1}^{n} w_i (t_i - \overline{t})}$$

with

$$\widehat{\theta}_1 = \overline{T}, \ \widehat{\theta}_2 = \widehat{\beta}_{\mathsf{IV}}.$$

We obtain

$$A(\theta) = \begin{pmatrix} 1 & 0 \\ 0 & \sigma_{X,T} \end{pmatrix}, \ B(\theta) = \begin{pmatrix} \sigma_T^2 & 0 \\ 0 & \sigma_T^2 \left(\sigma_{\varepsilon}^2 + \beta^2 \sigma_U^2 \right) \end{pmatrix}$$

• This yields the asymptotic covariance matrix

$$A(\theta)^{-1} B(\theta) \left(A(\theta)^{-1} \right)^{\mathsf{T}} = \begin{pmatrix} \sigma_{\mathcal{T}}^2 & 0 \\ 0 & \sigma_{\mathcal{T}}^2 \left(\sigma_{\varepsilon}^2 + \beta^2 \sigma_{U}^2 \right) / \sigma_{X,\mathcal{T}}^2 \end{pmatrix}$$

• When there is doubt about the magnitude of σ_U^2 , then we might want to estimate the joint asymptotic distribution of $\hat{\beta}_{IV}$ and $\hat{\beta}_{Y|W}$.

• *Example*. The sample *p*th quantile $\hat{\theta}_n = F_n^{-1}(p)$ satisfies

$$\psi\left(\mathsf{x}, heta
ight) = \mathbf{p} - \mathbb{I}\left(\mathsf{x} \leq heta
ight)$$

• We have
$$\sum_{i=1}^{n}\psi\left(x_{i}, heta
ight)=c_{n}=n\left(p-\mathcal{F}_{n}\left(\widehat{ heta}_{n}
ight)
ight)\leq1$$

 $\bullet\,$ This function is discontinuous at θ^* but we can have

$$A(\theta^{*}) = -\frac{\partial}{\partial \theta^{\mathsf{T}}} \mathbb{E} \left[\psi(X, \theta) \right] \Big|_{\theta = \theta^{*}} = -\frac{\partial}{\partial \theta^{\mathsf{T}}} \left[p - F(\theta) \right] \Big|_{\theta = \theta^{*}}$$
$$= f(\theta^{*}).$$

• We also have

$$B(\theta^*) = \mathbb{E}\left[p - \mathbb{I}\left(X \leq \theta^*\right)\right]^2 = p(1-p).$$

thus we have

$$V\left(heta^{*}
ight)=rac{p\left(1-p
ight)}{f\left(heta^{*}
ight)^{2}}.$$

- We could also stack any finite number of quantiles ψ functions together to get the joint asymptotic distribution of $(F_n^{-1}(p_1), \ldots, F_n^{-1}(p_k))$.
- However we cannot use $A(\widehat{\theta}_n, \mathbf{x})$ to estimate $A(\theta^*)$: in fact, the derivative of the *p*th quantile ψ function is zero everywhere except at the location of the jump discontinuity!
- To estimate f, we can use a kernel density estimator. An alternative consists of approximating ψ by a smooth ψ function.

 Example. The positive mean deviation from the median is defined to be

$$\widehat{\theta}_{1,n} = \frac{2}{n} \sum_{i=1}^{n} \left(x_i - \widehat{\theta}_{2,n} \right) \mathbb{I} \left(x_i \ge \widehat{\theta}_{2,n} \right)$$

where $\widehat{\theta}_{2,n}$ is the sample median.

• The ψ function is

$$\psi(x,\theta) = \left(\begin{array}{c} 2(x-\theta_2) \mathbb{I}(x \ge \theta_2) - \theta_1 \\ \frac{1}{2} - \mathbb{I}(x \le \theta_2) \end{array}\right).$$

 The 1st component of ψ is continuous everywhere but not differentiable at θ₂ = x. The 2nd component has a jump discontinuity at θ₂ = x. To get A(θ^{*}), we calculate

$$\mathbb{E}\left[\psi\left(X,\theta\right)\right] = \left(\begin{array}{c} 2\int_{\theta_2}^{\infty} \left(x-\theta_2\right)f\left(x\right)dx-\theta_1\\ \frac{1}{2}-F\left(\theta_2\right)\end{array}\right)$$

We write

$$2\int_{\theta_{2}}^{\infty} (x-\theta_{2}) f(x) dx - \theta_{1} = 2\int_{\theta_{2}}^{\infty} xf(x) dx - 2\theta_{2} \left[1 - F(\theta_{2})\right] - \theta_{1}.$$

The derivative of this expression with respect to θ_1 is -1, the derivative with θ_2 is

$$-2\theta_{2}f\left(\theta_{2}\right)-2\left[1-F\left(\theta_{2}\right)\right]+2\theta_{2}f\left(\theta_{2}\right).$$

It follows that

$$A\left(\theta^{*}\right) = \left(\begin{array}{cc}1 & 1\\0 & f\left(\theta_{2}^{*}\right)\end{array}\right), \ B\left(\theta^{*}\right) = \left(\begin{array}{cc}b_{11} & \frac{\theta_{1}^{*}}{2}\\\frac{\theta_{1}^{*}}{2} & \frac{1}{4}\end{array}\right)$$

where $b_{11} = 4 \int_{\theta_2}^{\infty} (x - \theta_2^*)^2 f(x) dx - \theta_1^{*2}$.

Finally we obtain

$$V\left(\theta^{*}\right) = \left(\begin{array}{cc} b_{11} - \frac{\theta_{1}^{*}}{f(\theta_{2}^{*})} + \frac{1}{4f(\theta_{2}^{*})^{2}} & \frac{\theta_{1}^{*}}{2f(\theta_{2}^{*})} - \frac{1}{4f(\theta_{2}^{*})^{2}} \\ \frac{\theta_{1}^{*}}{2f(\theta_{2}^{*})} - \frac{1}{4f(\theta_{2}^{*})^{2}} & \frac{1}{4f(\theta_{2}^{*})^{2}} \end{array}\right)$$