

Block 1  
Repetition from BSc courses  
LRM estimators & non-linear extensions  
Predictions from regression models

Advanced econometrics 1 4EK608  
Pokročilá ekonometrie 1 4EK416

Vysoká škola ekonomická v Praze

- ① Estimation methods, predictions from a model
  - Ordinary least squares
  - General properties of estimators
  - Method of moments
  - Maximum likelihood estimator
- ② Non-linear extensions to LRM, quantile regression
  - Non-linear regression models
  - Quantile regression
- ③ Predictions from a regression model
  - Predictions from a CLRM (repetition from BSc courses)
  - Predictions: general features,  $k$ FCV, Variance vs. Bias

# Linear regression model (LRM) and OLS estimation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

## LRM assumptions (for OLS estimation):

(Notation follows Greene, Econometric analysis, 7<sup>th</sup> ed.)

**A1 Linearity:**  $y_i = \beta_1 + \beta_2 x_{i2} + \cdots + \beta_K x_{iK} + \varepsilon_i$

LRM describes linear relationship between  $y_i$  and  $\mathbf{x}_i$ .

**A2 Full rank:** Matrix  $\mathbf{X}$  is an  $n \times K$  matrix with rank  $K$ .  
Columns of  $\mathbf{X}$  are linearly independent and  $n \geq K$ .

**A3 Exogeneity of regressors:**  $E[\varepsilon_i | \mathbf{X}] = 0$  (strict form).  
If relaxed to contemporaneous form in TS:  $E[\varepsilon_t | \mathbf{x}_t] = 0$ .  
Law of iterated expectations:  $E[\varepsilon_i | \mathbf{X}] = 0 \Rightarrow E[\varepsilon] = 0$ .

# Linear regression model (LRM) and OLS estimation

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

**LRM assumptions (continued):**

**A4 Homoscedastic & nonautocorrelated disturbances:**

$$E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'] = \sigma^2 \mathbf{I}_n$$

Homoscedasticity:  $\text{var}[\varepsilon_i | \mathbf{X}] = \sigma^2$ ,  $\forall i = 1, \dots, n$ .

Independent disturbances:  $\text{cov}[\varepsilon_t, \varepsilon_s | \mathbf{X}] = 0$ ,  $\forall t \neq s$ .

- GARCH models [i.e. ARCH(1):  $\text{var}[\varepsilon_t | \varepsilon_{t-1}] = \sigma^2 + \alpha \varepsilon_{t-1}^2$ ] do not violate the conditional variance assumption  $\text{var}[\varepsilon_i | \mathbf{X}] = \sigma^2$ . However,  $\text{var}[\varepsilon_t | \varepsilon_{t-1}] \neq \text{var}[\varepsilon_t]$ , with conditioning on  $\mathbf{X}$  omitted from notation but left as implicit.

**A5 DGP of  $\mathbf{X}$ :** Variables in  $\mathbf{X}$  may be fixed or random.

**A6 Normal distribution of disturbances:**

$$\boldsymbol{\varepsilon} | \mathbf{X} \sim N[\mathbf{0}, \sigma^2 \mathbf{I}_n].$$

# Ordinary least squares (OLS)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

The least squares estimator is unbiased (given A1 – A3):

$$\hat{\boldsymbol{\beta}} = \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon},$$

take expectations :

$$E[\mathbf{b}|\mathbf{X}] = \boldsymbol{\beta} + E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}|\mathbf{X}] = \boldsymbol{\beta}, \quad (\text{zero by A3}).$$

Variance of the least squares estimator (A1 – A4):

$$\text{var}[\mathbf{b}|\mathbf{X}] = \text{var}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}|\mathbf{X}]$$

because  $\text{var}(\boldsymbol{\beta}) = 0$ . Using A3 & A4:

$$= \mathbf{A}\boldsymbol{\sigma}^2\mathbf{I}_n\mathbf{A}' \quad \text{where } \mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

which is a matrix quadratic form for  $\text{var}(cZ) = c^2\text{var}(Z)$

$$= \boldsymbol{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$$

because  $(\mathbf{A}\mathbf{B})' = \mathbf{B}'\mathbf{A}'$ ; dim. compatible matrices  $\mathbf{A}, \mathbf{B}$ .

Normal distribution of the least squares estimator (A1 – A6):

$$\mathbf{b}|\mathbf{X} \sim N[\boldsymbol{\beta}, \boldsymbol{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}].$$

## Estimators and estimation methods:

- LRM is not the only type of regression model.
- OLS is not the only useful estimator.
- Let's approach estimators and their properties more generally.

(again, notation follows Greene, Econometric analysis.)

# Estimators and estimation methods

Notation/definitions:

- $\mathbf{x}'_j = (x_{1j}, \dots, x_{nj})$  - random sample of  $n$  observations.
- $\boldsymbol{\theta}$  - population parameter [unknown parameter(s)]
- $f(\mathbf{x}_j, \boldsymbol{\theta})$ : probability distribution function
- $\hat{\boldsymbol{\theta}}$  is some estimator of  $\boldsymbol{\theta}$

Basic notions:

- All estimators have sampling distributions  
mean:  $E(\hat{\theta})$   
variance:  $E[(\hat{\theta} - E(\hat{\theta}))^2]$ , etc.
- Estimators  $\times$  estimate
- Generally, many estimators exist for a given parameter.  
Population mean example (two sample-based estimators):

$$\hat{\theta}_1 = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\hat{\theta}_2 = \tilde{x} = \frac{1}{2}(x_{max} + x_{min})$$

# Properties of estimators - classification:

- **Unbiasedness:** can be described as  $E(\hat{\theta}) = \theta$ .  
Occasionally useful – in finite (small) sample context.  
**Asymptotic unbiasedness** (large sample property):  
not very useful, discussion would be directed towards consistency (which is a far more desirable feature).
- **Consistency:**  $\text{plim}(\hat{\theta}) = \theta$ .  
 $\hat{\theta} \rightarrow \theta$  as  $n \rightarrow \infty$ : vector  $\hat{\theta}$  is at least asymptotically unbiased and  $\text{plim}(\text{var}(\hat{\theta})) = \mathbf{0}$  [i.e.  $\text{var}(\hat{\theta}) \rightarrow \mathbf{0}$  as  $n \rightarrow \infty$ ].
  - Consistent estimators: unbiased or asymptotically unbiased & their variance shrinks to zero as sample size grows.
  - Minimal requirement for estimator used in statistics or econometrics.
  - If some estimator is not consistent, then it does not provide relevant estimates of population  $\theta$  values, even with unlimited data, i.e. as  $n \rightarrow \infty$ .
  - Unbiased estimators are not necessarily consistent.
  - Biased and consistent estimators are often useful (small-sample bias, yet consistent: IVR, ML, etc.).



# Properties of estimators - classification:

- **Efficiency:** an estimator is efficient if it is unbiased and no other unbiased estimator has a smaller variance. Often difficult to prove, we usually simplify the concept to *relative efficiency* (e.g.: efficiency with respect to linear unbiased estimators, etc.).

**Asymptotic efficiency:** holds for an estimator that is asymptotically unbiased and no other asymptotically unbiased estimator has smaller asymptotic variance.

- **Normality, asymptotic normality:** basis for most statistical inference performed with common estimators.

# Estimators and estimation methods

**Extremum estimator:** obtained as the optimizer of some criterion function  $q(\boldsymbol{\theta}|\mathbf{data})$ . Most common estimators:

$$\text{LS } \hat{\boldsymbol{\theta}}_{LS} = \operatorname{argmax} \left[ -\frac{1}{n} \sum_{i=1}^n (y_i - h(\mathbf{x}_i, \boldsymbol{\theta}_{LS}))^2 \right],$$

$$\text{ML } \hat{\boldsymbol{\theta}}_{ML} = \operatorname{argmax} \left[ \frac{1}{n} \sum_{i=1}^n \log f(y_i | \mathbf{x}_i, \boldsymbol{\theta}_{ML}) \right],$$

$$\text{GMM } \hat{\boldsymbol{\theta}}_{GMM} = \operatorname{argmax} \left[ -\overline{\mathbf{m}}(\mathbf{data}, \boldsymbol{\theta}_{GMM})' \mathbf{W} \overline{\mathbf{m}}(\mathbf{data}, \boldsymbol{\theta}_{GMM}) \right],$$

where  $h(\cdot)$  is a function (linear/non-linear  $\rightarrow$  OLS/NLS),

$f(\cdot)$  is a probability density function (pdf),

$\overline{\mathbf{m}}$  denotes sample moments,

$\mathbf{W}$  is a convenient positive definite matrix.

LS and ML estimators belong to a class of **M estimators** (type of extremum estimators where objective function is a sample average).

Assumptions for asymptotic properties of extremum estimators:

- 1 **Parameter space:** must be convex and the parameter vector that is the object of estimation must be point in its interior. Gaps and nonconvexities in parameter spaces would generally collide with estimation algorithms (settings such as  $\sigma^2 \geq 0$  are OK).
- 2 **Criterion function:** must be concave in the parameters (concave in the neighborhood of the true parameter vector). Criterion functions need not be globally concave. In such situation, there may be multiple local optima (often associated with poor model specification).

Assumptions for asymptotic properties of extremum estimators:

**3 Identifiability of parameters:** has a relatively complex technical definition (anything like “true parameters  $\theta_0$  are identified if...” is problematic - leads to a paradox if condition is not met). Simple way to secure identification:

- **LS:** for a given set of any two different parameter vectors  $\theta$  and  $\theta_0$ , a vector of observations  $\mathbf{x}_i$  must exist (for some  $i$ ), leading to different conditional mean function ( $\hat{y}_i$ ).
- **ML:** For any two parameter vectors  $\theta \neq \theta_0$ , a data vector  $(y_i, \mathbf{x}_i)$  must exist, which generates different values of density function:  $f(y_i|\mathbf{x}_i, \theta) \neq f(y_i|\mathbf{x}_i, \theta_0)$ .

Note: identifiability does not rule out possibility of:  
 $f(y_i|\mathbf{x}_i, \theta) = f(y_\ell|\mathbf{x}_\ell, \theta)$ , where,  $y_i = y_\ell$ ,  $\mathbf{x}_i \neq \mathbf{x}_\ell$ .

- **GMM:** sufficient condition for identification:  
 $E[\overline{\mathbf{m}}(\mathbf{data}, \theta)] \neq \mathbf{0}$  if  $\theta \neq \theta_0$ .

Assumptions for asymptotic properties of extremum estimators:

## 4 Behavior of the data: Grenander conditions for well-behaved data:

- G1 For each  $\mathbf{x}_k$  column of  $\mathbf{X}$  and  $d_{nk}^2 = \mathbf{x}_k' \mathbf{x}_k = \sum_{i=1}^n x_{ik}^2$ , it must hold that:  $\lim_{n \rightarrow \infty} d_{nk}^2 = +\infty$ .  
Sum of squares continue to grow with sample size, i.e.  $\mathbf{x}_k$  does not degenerate into a series of 0.
- G2 The  $\lim_{n \rightarrow \infty} x_{ik}^2 / d_{nk}^2 = 0$  for all  $i = 1, 2, \dots, n$ . Single observations become less important as sample size grows.  
No single observation will dominate  $\mathbf{x}_k' \mathbf{x}_k$ .
- G3 Let  $\mathbf{C}_n$  be sample correlation matrix of the columns in  $\mathbf{X}$  (excluding the intercept, if present). Then  $\lim_{n \rightarrow \infty} \mathbf{C}_n = \mathbf{C}$  where  $\mathbf{C}$  is positive definite. This implies that the full rank condition for  $\mathbf{X}$  (A2) is not asymptotically violated.

## Quick convergence recap (terminology):

- Convergence in probability: a sequence of random variables  $X_1, X_2, X_3, \dots$  converges in probability to a random variable  $X$ , denoted as  $X_n \xrightarrow{p} X$  [or  $\text{plim}(X_n) = X$ ], if:

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0, \quad \forall \epsilon > 0.$$

- Convergence in distribution: a weaker type of convergence. It states that the CDF of  $X_n$  converges to the CDF of  $X$  as  $n$  goes to infinity (does not require dependency between  $X_n$  and  $X$ ).  $X_n \xrightarrow{d} X$ , if:

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x), \quad F_X(x) \text{ continuous.}$$

## Theorem: Consistency of M estimators

If:

- (a) the parameter space is convex and the true parameter vector is a point in its interior,
- (b) the criterion function is concave,
- (c) the parameters are identified by the criterion function,
- (d) the data are well behaved,

then the M estimator converges in probability to the true parameter vector.

# Estimators and estimation methods

## Theorem: Asymptotic normality of M estimators

If:

- (a)  $\hat{\boldsymbol{\theta}}$  is a consistent estimator of  $\boldsymbol{\theta}_0$  where  $\boldsymbol{\theta}_0$  is a point in the interior of the parameter space  $\boldsymbol{\Theta}$ ,
- (b)  $q(\boldsymbol{\theta}|\mathbf{data})$  is concave and twice continuously differentiable in  $\boldsymbol{\theta}$  in a neighborhood of  $\boldsymbol{\theta}_0$ ,
- (c)  $\sqrt{n} [\partial q(\boldsymbol{\theta}_0|\mathbf{data})/\partial \boldsymbol{\theta}_0] \xrightarrow{d} N(\mathbf{0}, \boldsymbol{\Phi})$ ,
- (d)  $\lim_{n \rightarrow \infty} \Pr [ |(\partial^2 q(\boldsymbol{\theta}|\mathbf{data})/\partial \theta_k \partial \theta_m) - h_{km}(\boldsymbol{\theta})| > \varepsilon ] = 0 \ \forall \varepsilon > 0$  for any  $\boldsymbol{\theta}$  in  $\boldsymbol{\Theta}$ ;  $h_{km}(\boldsymbol{\theta})$  is a continuous finite valued function of  $\boldsymbol{\theta}$ ,
- (e) the matrix of elements  $\mathbf{H}(\boldsymbol{\theta})$  is nonsingular at  $\boldsymbol{\theta}_0$ ,

then  $\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \xrightarrow{d} N\{\mathbf{0}, [\mathbf{H}^{-1}(\boldsymbol{\theta}_0)\boldsymbol{\Phi}\mathbf{H}^{-1}(\boldsymbol{\theta}_0)]\}$ .

where  $\boldsymbol{\Phi}$  is a variance-covariance matrix,

and  $\mathbf{H}(\boldsymbol{\theta}_0) = \partial^2 q(\boldsymbol{\theta}|\mathbf{data})/\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'$  is a Hessian (evaluated at  $\boldsymbol{\theta}_0$ ).



- Method of moments (MM)
- Generalized method of moments (GMM)

# Method of moments

- With the method of moments, we simply estimate population moments by corresponding sample moments.
- Under very general conditions, sample moments are consistent estimators of the corresponding population moments, but NOT necessarily unbiased estimators.

## Application example 1

Sample covariance is a consistent estimator of population covariance.

## Application example 2

OLS estimators we have used for parameters in the CLRM can be derived by the method of moments.

# Method of moments

## Method of moments (MM)

Population moments for a stochastic variable  $X$

- $E(X^r)$ :  $r^{th}$  population moment about zero
- $E(X)$ : population mean: 1<sup>st</sup> population moment about zero
- $E[(X - E(X))^2]$ : population variance is the second moment about the mean

Sample moments for sample observations  $(x_1, x_2, \dots, x_n)$

- $\frac{\sum_{i=1}^n x_i^r}{n}$ :  $r^{th}$  sample moment about zero
- $\frac{\sum_{i=1}^n x_i}{n} = \bar{x}$ : sample mean is the first moment about zero
- $\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$ : sample variance is the second sample moment about the mean

- For MM, the usual linear model assumption (concerning 1<sup>st</sup> population moment)  $E[\mathbf{x}_i \varepsilon_i] = \mathbf{0}$  implies:

$$E[\mathbf{x}_i(y_i - \mathbf{x}_i' \boldsymbol{\beta})] = \mathbf{0},$$

which constitutes a **population moment equation**:

$$E[\mathbf{x}_i(y_i - \mathbf{x}_i' \boldsymbol{\beta})] = E[\mathbf{m}(\boldsymbol{\beta})] = \mathbf{0},$$

and the corresponding sample (empirical) moment equation can be formalized as:

$$\left[ \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i(y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}) \right] = \overline{\mathbf{m}}(\hat{\boldsymbol{\beta}}) = \mathbf{0}.$$

# Method of moments

For a LRM with  $K$  exogenous regressors, MM sample equations can be cast as:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK}) = 0$$

$$\frac{1}{n} \sum_{i=1}^n x_{i2} (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK}) = 0$$

...

$$\frac{1}{n} \sum_{i=1}^n x_{iK} (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK}) = 0$$

- Removing  $\frac{1}{n}$  elements from equations does not affect the solution.
- This is a system of  $K$  equations with  $K$  unknown parameters  $\beta_j$ .
- The set of moment equations is equivalent to 1<sup>st</sup> order conditions for the OLS estimator:

$$\min_{\hat{\beta}} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK})^2$$

# Generalized method of moments

- GMM is a very general class of estimators, includes many other estimators as a special case (IVR, simultaneous equations, Arellano-Bond estimator for dynamic panels).
- For single equation linear models, GMM may be conveniently described using the instrumental variable case:

For the LRM  $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$ ,

- We abandon the assumption  $E[\mathbf{x}_i \varepsilon_i] = 0$  and
- we replace it by  $E[\mathbf{z}_i \varepsilon_i] = 0$ .
- Hence, columns of  $\mathbf{X}$  ( $n \times K$ ) are potentially endogenous and  $\mathbf{Z}$  ( $n \times L$ ) is a matrix of exogenous instruments.

All exogenous columns in  $\mathbf{X}$  are repeated in  $\mathbf{Z}$  and each endogenous column in  $\mathbf{X}$  is replaced in  $\mathbf{Z}$  by at least one instrument (exogenous variable not present in  $\mathbf{X}$ ).

# Generalized method of moments

- GMM equations can be cast by analogy to the MM case:

we start by  $E[\mathbf{z}_i \varepsilon_i] = \mathbf{0}$ , which implies a **population moment equation** (matrix form):

$$E[\mathbf{z}_i(y_i - \mathbf{x}_i' \boldsymbol{\beta})] = E[\mathbf{m}(\boldsymbol{\beta})] = \mathbf{0},$$

and corresponding sample (empirical) moment equation:

$$\left[ \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i(y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}) \right] = \overline{\mathbf{m}}(\hat{\boldsymbol{\beta}}) = \mathbf{0}.$$

# Generalized method of moments

GMM empirical equations can also be cast as:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK}) = 0$$

$$\frac{1}{n} \sum_{i=1}^n z_{i2} (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK}) = 0$$

...

$$\frac{1}{n} \sum_{i=1}^n z_{iL} (y_i - \hat{\beta}_1 - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_K x_{iK}) = 0$$

- First column of  $\mathbf{Z}$  is assumed to be a vector of ones (same as for  $\mathbf{X}$ ).
- For  $\mathbf{Z} = \mathbf{X}$  as a special case, the above equations are identical to MM (shown previously) and the solution is identical to the OLS estimator:  $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ .
- For  $\mathbf{Z} \neq \mathbf{X}$ , where  $\mathbf{Z}$  is  $(n \times L)$  and  $\mathbf{X}$  is  $(n \times K)$ , three identification possibilities have to be considered.



## Identification of GMM equations

1 **Underidentified:** with  $L < K$ , there are fewer moment equations than unknown parameters ( $\beta_j$ ). Without additional information (parameter restrictions), there is no solution to the system of GMM equations.

2 **Exactly identified:** for  $L = K$ , single solution exists:

$$\left[ \frac{1}{n} \sum_{i=1}^n z_i (y_i - \mathbf{x}_i' \hat{\beta}) \right] = \overline{\mathbf{m}}(\hat{\beta}) = \mathbf{0},$$

can be conveniently re-written as:

$$\overline{\mathbf{m}}(\hat{\beta}) = \left( \frac{1}{n} \mathbf{Z}' \mathbf{y} \right) - \left( \frac{1}{n} \mathbf{Z}' \mathbf{X} \right) \hat{\beta} = \mathbf{0}$$

and the solution yields the familiar IV estimator:

$$\hat{\beta} = (\mathbf{Z}' \mathbf{X})^{-1} \mathbf{Z}' \mathbf{y}.$$

## Identification of GMM equations (continued)

- 3 With  $L > K$ , there is no unique solution to the equation system  $\overline{\mathbf{m}}(\hat{\beta}) = \mathbf{0}$ .

One intuitive solution is the “least squares approach”:

$$\min_{\beta} \left( \overline{\mathbf{m}}(\hat{\beta})' \overline{\mathbf{m}}(\hat{\beta}) \right)$$

Through the first order conditions, we obtain a GMM estimator as

$$\hat{\beta} = [(\mathbf{X}'\mathbf{Z})(\mathbf{Z}'\mathbf{X})]^{-1} (\mathbf{X}'\mathbf{Z})\mathbf{Z}'\mathbf{y}.$$

# Generalized method of moments

## GMM - consistency conditions

- **Convergence of the moments:** Empirical (sample) moments converge in probability to their population counterparts. DGP meets the conditions for LLN.

$$\bar{\mathbf{m}}(\boldsymbol{\beta}) = \frac{1}{n} (\mathbf{Z}'\mathbf{y} - \mathbf{Z}'\mathbf{X}\boldsymbol{\beta}) \xrightarrow{p} \mathbf{0}.$$

- **Identification:** For any  $n \geq K$  and  $\boldsymbol{\beta}_1 \neq \boldsymbol{\beta}_2$  it holds that  $\bar{\mathbf{m}}(\boldsymbol{\beta}_1) \neq \bar{\mathbf{m}}(\boldsymbol{\beta}_2)$ . Three implications:
  - **Order condition:**  $L \geq K$ . Number of moment equations at least as large as number of parameters.
  - **Rank condition:** matrix  $\mathbf{G}(\boldsymbol{\beta}) = \partial \bar{\mathbf{m}}(\boldsymbol{\beta}) / \partial \boldsymbol{\beta}'$  (i.e.  $\frac{1}{n} \mathbf{Z}'\mathbf{X}$ ) is a  $L \times K$  matrix with row rank equal to  $K$ .
  - **Uniqueness:** unique solution/optimizer exists.
- **Limiting Normal distribution for the sample moments:** Population moments obey central limit theorem (CLT) or some similar variant.

## GMM - final remarks & summary

- GMM-based asymptotic covariance matrix of  $\hat{\beta}$  is discussed in Greene (Econometric analysis, ch. 13.6) for the classical, heteroscedastic and generalized case (includes TS-based estimation).
- GMM is robust to differences in “specification” of the data generating process (DGP).  $\rightarrow$  i.e. sample mean or sample variance estimate their population counterparts (assuming they exist) regardless of DGP.
- GMM is free from distributional assumptions. “Cost” of this approach: if we know the specific distribution of a DGP, GMM does not make use of such information  $\rightarrow$  inefficient estimates.
- Alternative approach: method of maximum likelihood utilizes distributional information and is more efficient (provided this information is available & valid).

# Maximum likelihood estimator

- Maximum likelihood estimator (MLE)
- Normal distribution & MLE

# Maximum likelihood estimator

## Maximum likelihood estimator – single parameter

For a stochastic variable  $y$  with a known distribution, described by a single  $\theta$  parameter:

- $f(y|\theta)$  is the pdf/pmf of  $y$ , conditioned on parameter  $\theta$ .  
(pmf: probability mass function, discrete probability density f.)
- For  $n$  *iid* observations, joint density of this process:

$$f(y_1, y_2, \dots, y_n | \theta) = \prod_{i=1}^n f(y_i | \theta) = L(\theta | \mathbf{y})$$

is the likelihood function.

- We estimate  $\theta$  by maximizing  $L(\theta | \mathbf{y})$  with respect to the parameter (1<sup>st</sup> order conditions). Solution (MLE) often denoted as  $\hat{\theta}_{\text{ML}}$ .
- For maximization (MLE), it is usually simpler to work with a log-transformed likelihood function:

$$\log L(\theta | \mathbf{y}) = \sum_{i=1}^n \log f(y_i | \theta).$$

## MLE – Poisson distribution example

- Consider 10 *iid* observations from a Poisson distribution:  
 $\mathbf{y}' = (5, 0, 1, 1, 0, 3, 2, 3, 4, 1)$ .
- The pmf:  $f(y_i|\lambda) = \Pr(Y = y_i) = \frac{e^{-\lambda}\lambda^{y_i}}{y_i!}$ .
- Likelihood function:  $L(\lambda|\mathbf{y}) = \prod_{i=1}^n \frac{e^{-\lambda}\lambda^{y_i}}{y_i!} = \frac{e^{-10\lambda}\lambda^{\sum_{i=1}^{10} y_i}}{\prod_{i=1}^{10} y_i!}$ .
- $\log L$ :  $\log L(\lambda|\mathbf{y}) = -n\lambda + \log \lambda \sum_{i=1}^n y_i - \sum_{i=1}^n \log(y_i!),$
- 1<sup>st</sup> order condition:  $\frac{\partial \log L(\lambda|\mathbf{y})}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n y_i = 0.$
- From 1<sup>st</sup> order condition:  $\hat{\lambda}_{\text{ML}} = \bar{y}_n.$   
For our empirical example:  $\hat{\lambda}_{\text{ML}} = 2.$

# Maximum likelihood estimator

## Maximum likelihood estimator – vector of parameters

- $\theta = (\theta_1, \dots, \theta_m)'$
- $L = L(\theta_1, \theta_2, \dots, \theta_m | y_1, y_2, \dots, y_n)$
- We find MLEs of the  $m$  parameters by partially differentiating the likelihood function  $L$  (often,  $\log L$  is used) with respect to each  $\theta$  and then setting all the partial derivatives obtained to zero.



## LRM parameters & Normal distribution

- $L(\boldsymbol{\theta}|\mathbf{data}) = L(\boldsymbol{\beta}, \sigma^2|y_i, \mathbf{x}_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{x}_i'\boldsymbol{\beta})^2}{2\sigma^2}}$
- In matrix form, the log likelihood function is:

$$LL(\boldsymbol{\beta}, \sigma^2|\mathbf{y}, \mathbf{X}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Recall that:

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}'\mathbf{y} - 2\boldsymbol{\beta}'\mathbf{X}'\mathbf{y} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}$$

and

$$\frac{\partial(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{\partial\boldsymbol{\beta}'} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\boldsymbol{\beta}.$$

## LRM parameters & Normal distribution (continued)

$$LL(\beta, \sigma^2 | \mathbf{y}, \mathbf{X}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$$

- 1<sup>st</sup> order conditions:

- $\frac{\partial LL}{\partial \beta'} = \frac{1}{2\sigma^2} [2\mathbf{X}'\mathbf{y} - 2\mathbf{X}'\mathbf{X}\beta] = \mathbf{0}$

is solved by:

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$$

- $\frac{\partial LL}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} [(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)] = 0$

is solved by:

$$\hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})}{n} = \frac{\mathbf{u}'\mathbf{u}}{n} = \frac{\text{SSR}}{n}.$$

Note: the MLE estimate  $\hat{\sigma}^2$  is biased downwards in small samples, as the unbiased estimate is equal to  $\text{SSR}/(n - K)$ .

## MLE assumptions (quick recap)

- **Parameter space:** Gaps and nonconvexities in parameter spaces would generally collide with estimation algorithms.
- **Identifiability:** The parameter vector  $\theta$  is identified (estimable), if for two vectors,  $\theta^* \neq \theta$  and for some data observations  $\mathbf{x}$ ,  $L(\theta^*|\mathbf{x}) \neq L(\theta|\mathbf{x})$ .
- **Well-behaved data:** Laws of large numbers (LLN) apply. Some form of CLT can be applied to the gradient (i.e. for the estimation method).
- **Regularity conditions:** “well behaved” derivatives of  $f(y_i|\theta)$  with respect to  $\theta$  (see Greene, chapter 14.4.1).

# Maximum likelihood estimator

**MLE properties** (if assumptions are met)

- **Consistency:**  $\text{plim}(\hat{\theta}) = \theta_0$  ( $\theta_0$  is the true parameter)
- **Asymptotic normality** of  $\hat{\theta}$
- **Asymptotic efficiency:**  $\hat{\theta}$  is asymptotically efficient and achieves the Cramér-Rao lower bound for consistent estimators (see Greene, chapter 14.4.5)
- **Invariance:** MLE of  $\gamma_0 = c(\theta_0)$  is  $c(\hat{\theta})$  if  $c(\theta_0)$  is a continuous and countinuously differentiable function.  
(empirical advantages: we can use reparameterization in MLE, e.g.  $\gamma_j = 1/\theta_j$  or  $\theta^2 = 1/\sigma^2$ ).

# Maximum likelihood estimator

**MLE properties** (Normal distribution,  $\text{var}(\hat{\boldsymbol{\theta}})$ ):

Under the above assumptions, asymptotic variance-covariance matrix of  $\hat{\boldsymbol{\theta}}$  is the inverse of the Information matrix:

$$\text{var}(\hat{\boldsymbol{\theta}}) = \mathbf{I}[\hat{\boldsymbol{\theta}}]^{-1} = \begin{bmatrix} \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{2\hat{\sigma}^4}{n} \end{bmatrix},$$

where  $\mathbf{I}[\hat{\boldsymbol{\theta}}] = -[\mathbf{H}(\hat{\boldsymbol{\theta}})]$ . MLE gives the familiar formula for  $\text{var}(\hat{\beta}) = \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$ , and a simple expression for the variance of  $\hat{\sigma}^2$ . The two zero vectors come from LRM exogeneity assumptions. (see Greene, 7<sup>th</sup> ed., ch. 14.4)

- The square roots of diagonal elements of  $\mathbf{I}[\hat{\boldsymbol{\theta}}]^{-1}$  give estimates of the standard errors of parameter estimates.
- We can construct simple  $z$ -scores to test the null hypothesis concerning any individual parameter, just as in OLS, but using the normal instead of the  $t$ -distribution.

## MLE - inference, three classic tests:

Consider MLE of parameter  $\boldsymbol{\theta}$  and a test of the hypothesis  $H_0 : \boldsymbol{h}(\boldsymbol{\theta}) = \mathbf{0}$ . Recall that ML parameter estimates are asymptotically normally distributed.

- 1 **Likelihood ratio test:** If the restriction  $\boldsymbol{h}(\boldsymbol{\theta}) = \mathbf{0}$  is valid, then imposing it should not lead to a large reduction in the log-likelihood function.

$$LR = 2(LL_U - LL_R) \underset{H_0}{\sim} \chi^2(r),$$

where  $LL_U$  is the  $LL$  of unconstrained model,  $LL_R$  denotes restricted model and  $r$  is the number of restrictions imposed. To do this test you have to estimate two models (one nested) and get the results of both.

## MLE - inference, three classic tests:

We have an unrestricted ML estimate  $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_m)'$ ,  
and test of the hypothesis  $H_0 : \mathbf{h}(\boldsymbol{\theta}) = \mathbf{q}$ ,  
where  $\mathbf{q}$  is a  $(r \times 1)$  vector function of  $\boldsymbol{\theta}$  (linear/non-linear  
restrictions, continuous partial derivatives assumed).

- 2 Wald test:** If restriction  $\mathbf{h}(\boldsymbol{\theta}) = \mathbf{q}$  is valid, then  $\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{q}$   
should be close to zero since MLE is consistent.

$$W = [\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{q}]' \left[ \text{Asy.Var}[\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{q}] \right]^{-1} [\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{q}] \underset{H_0}{\sim} \chi^2(r),$$

where the estimated

$$\text{Asy.Var}[\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{q}] = \left[ \frac{\partial \mathbf{h}(\hat{\boldsymbol{\theta}})}{\partial \hat{\boldsymbol{\theta}}'} \right] \text{Est.Asy.Var}(\hat{\boldsymbol{\theta}}) \left[ \frac{\partial \mathbf{h}(\hat{\boldsymbol{\theta}})}{\partial \hat{\boldsymbol{\theta}}'} \right]'$$

# Maximum likelihood estimator

## MLE - inference, three classic tests:

We have a ML estimate  $\hat{\boldsymbol{\theta}}_R$  – i.e. ML estimation of the restricted model, under  $H_0 : \mathbf{h}(\boldsymbol{\theta}) = \mathbf{0}$ ,

- 3 Lagrange multiplier test:** If the restriction is valid, then the restricted estimator should be near the point that maximizes the log-likelihood. Therefore, the slope of the log-likelihood function should be near zero at the restricted estimator. The test is based on the slope of the log-likelihood at the point where the function is maximized subject to the restriction.

$$LM = \left( \frac{\partial \log L(\hat{\boldsymbol{\theta}}_R)}{\partial \hat{\boldsymbol{\theta}}_R} \right)' \mathbf{I}[\hat{\boldsymbol{\theta}}_R]^{-1} \left( \frac{\partial \log L(\hat{\boldsymbol{\theta}}_R)}{\partial \hat{\boldsymbol{\theta}}_R} \right) \underset{H_0}{\sim} \chi^2(r),$$

where  $-\mathbf{I}[\hat{\boldsymbol{\theta}}_R] = \partial^2 LL(\boldsymbol{\theta}) / \partial \boldsymbol{\theta}' \partial \boldsymbol{\theta}$  evaluated at  $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_R$ .



## MLE - inference, three classic tests:

- The  $\chi^2$  distributions of the three test statistics are asymptotically valid.
- The three tests are asymptotically equivalent, but may differ in small samples:
- $W \geq LR \geq LM$ .
- Hence, in finite samples,  $LR$  rejects  $H_0$  less often than  $W$  but more often than  $LM$ .
- The above tests are discussed in ML context, i.e. with a known distribution of the variable/error term (ML parameter estimates are asymptotically normally distributed).

## MLE – summary

- MLE is only possible if we know the form of the probability distribution function for the population (Normal, Poisson, Negative Binomial, etc.).
- MLE has the large sample properties of consistency and asymptotic efficiency. There is no guarantee of desirable small-sample properties.
- Under CLRM assumptions (A1 – A6), ML estimator is identical to OLS estimator (for  $\hat{\beta}$ ).

# Non-linear extensions to LRM, quantile regression

- Non-linear regression models
- Quantile regression

## Nonlinear regression model:

$$y_i = h(\mathbf{x}_i, \boldsymbol{\beta}) + \varepsilon_i$$

- Linear model is a special case of the nonlinear model.
  - $y_i = h(\mathbf{x}_i, \boldsymbol{\beta}) + \varepsilon_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$ .
  - Linear models: linear in parameters. Definition includes non-linear regressors such as  $x_i^2$ , etc.
  - Many nonlinear models can be transformed into linear models (log-transformation)
- For nonlinear models that cannot be transformed into LRM, nonlinear LS (NLS) are available.
- $\partial h(\mathbf{x}_i, \boldsymbol{\beta}) / \partial \mathbf{x}$  is no longer equal to  $\boldsymbol{\beta}$  (interpretation based on estimated model ...)

## Assumptions relevant to the nonlinear regression model

- 1 **Functional form:** The conditional mean function for  $y_i$ , given  $\mathbf{x}_i$  is:

$$\mathbf{E}[y_i|\mathbf{x}_i] = h(\mathbf{x}_i, \boldsymbol{\beta}) , \quad i = 1, 2, \dots, n$$

- 2 **Identifiability of model parameters:** The parameter vector in the model is identified (estimable) if there is no nonzero parameter  $\boldsymbol{\beta}_0 \neq \boldsymbol{\beta}$  such that  $h(\mathbf{x}_i, \boldsymbol{\beta}_0) = h(\mathbf{x}_i, \boldsymbol{\beta})$  for all  $\mathbf{x}_i$ .
- 3 **Zero mean of the disturbance:** For  $y_i = h(\mathbf{x}_i, \boldsymbol{\beta}) + \varepsilon_i$ , we assume

$$\mathbf{E}[\varepsilon_i|h(\mathbf{x}_i, \boldsymbol{\beta})] = 0 , \quad i = 1, 2, \dots, n$$

i.e. disturbance at observation  $i$  is uncorrelated with the conditional mean function.

## Assumptions relevant to the nonlinear regression model

### 4 Homoscedasticity and non-autocorrelation:

conditional homoscedasticity:

$$\mathbf{E}[\varepsilon_i^2 | h(\mathbf{x}_i, \boldsymbol{\beta})] = \sigma^2, \quad i = 1, 2, \dots, n$$

non-autocorrelation:

$$\mathbf{E}[\varepsilon_t \varepsilon_s | h(\mathbf{x}_t, \boldsymbol{\beta}), h(\mathbf{x}_s, \boldsymbol{\beta})] = 0, \quad \text{for all } t \neq s$$

## Assumptions relevant to the nonlinear regression model

- 5 **Data generating process:** DGP for  $\mathbf{x}_i$  is assumed to be a well-behaved population such that first and second sample moments of the data can be assumed to converge to fixed, finite population counterparts. The crucial assumption is that the process generating  $\mathbf{x}_i$  is strictly exogenous to that generating  $\varepsilon_i$
- 6 **Underlying probability model** There is a well-defined probability distribution generating  $\varepsilon_i$ . At this point, we assume only that this process produces a sample of uncorrelated, identically (marginally) distributed random variables  $\varepsilon_i$  with mean zero and variance  $\sigma^2$  conditioned on  $h(\mathbf{x}_i, \boldsymbol{\beta})$ . Hence, our statement of the model is **semi-parametric** (i.e. specific distributional assumption on residuals are replaced by weaker assumptions).

## NLS: estimator of the nonlinear regression model

- NLS:      min:     $S(\beta) = \sum [y_i - h(\mathbf{x}_i, \beta)]^2$
- Using the standard procedure, we can get  $k$  first order conditions for the minimization:

$$\frac{\partial S(\beta)}{\partial \beta} = 2 \cdot \sum_{i=1}^n [y_i - h(\mathbf{x}_i, \beta)] \frac{\partial h(\mathbf{x}_i, \beta)}{\partial \beta} = \mathbf{0}$$

- The above first order conditions are also moment conditions and this defines the NLS estimator as a GMM estimator.



## **NLS: estimator of the nonlinear regression model**

- NLS being a GMM estimator allows us to deduce that the NLS estimator has good large sample properties: consistency and asymptotic normality (if assumptions are fulfilled).
- Hypothesis testing: The principal testing procedure is the Wald test, which relies on the consistency and asymptotic normality of the estimator. Likelihood ratio and LM tests can also be constructed.

# Nonlinear regression: computing NLS estimates

For nonlinear models, a closed-form solution (NLS estimator) usually does not exist.

- Most of the nonlinear maximization problems are solved by an **iterative algorithm**.
- The most commonly used of iterative algorithms are **gradient methods**.
- The template for most gradient methods in common use is the **Newton's method**.
- Look at your software packages which methods are available for computing NLS estimates.

# Nonlinear regression: examples

- LRM on TS with autocorrelation:

$$y_t = \mathbf{x}'_t \boldsymbol{\beta} + u_t,$$

$$u_t = \rho u_{t-1} + \varepsilon_t,$$

$$y_t = \mathbf{x}'_t \boldsymbol{\beta} + \rho u_{t-1} + \varepsilon_t \quad \text{note: } u_{t-1} = y_{t-1} - \mathbf{x}'_{t-1} \boldsymbol{\beta},$$

hence:

$$y_t = \rho y_{t-1} + \mathbf{x}'_t \boldsymbol{\beta} + \rho(\mathbf{x}'_{t-1} \boldsymbol{\beta}) + \varepsilon_t,$$

which is non-linear in parameters  $(\rho \boldsymbol{\beta})$ .

- Non-linear consumption function example:

$$cons_i = \beta_1 + \beta_2 inc_i^{\beta_3} + \varepsilon_i$$

special case: model is linear for  $\beta_3 = 1$   
(such assumption can be tested).

# Nonlinear regression: examples

Examples 7.4 & 7.8 (Greene):

Analysis of a Nonlinear Consumption Function

OLS version: for  $\beta_3 = 1$ .

Dependent Variable: REALCONS				
Method: Least Squares (Marquard - EViews legacy)				
Date: 09/19/16 Time 16:31				
Sample 1950Q1 2000Q4				
Included observations: 204				
REALCONS=C(1)+C(2)*REALDPI				
	Coefficient	Std.Error	t-Statistic	Prob.
C(1)	-80.35475	14.30585	-5.616915	0.0000
C(2)	0.921686	0.003872	238.0540	0.0000
R-squared	0.996448	Mean dependent var		2999.436
Adjusted R-squared	0.996431	S.D. dependent var		1459.707
S.E. of regression	87.20983	Akaike info criterion		11.78427
Sum squared resid	1536322	Schwarz criterion		11.81680
Log likelihood	-1199.995	Hannan-Quinn criter.		11.79743
F-statistics	56669.72	Durbin-Watson stat		0.092048
Prob(F-statistics)	0.000000			

# Nonlinear regression: examples

Examples 7.4 & 7.8 (Greene):

Analysis of a Nonlinear Consumption Function

NLS with starting values equal to 0

Dependent Variable: REALCONS				
Method: Least Squares (Marquard - EViews legacy)				
Sample 1950Q1 2000Q4    Included observations: 204				
Convergence achieved after 200 iterations				
REALCONS=C(1)+C(2)*REALDPI^C(3)				
	Coefficient	Std.Error	t-Statistic	Prob.
C(1)	458.7991	22.50140	20.38980	0.0000
C(2)	0.100852	0.010910	9.243667	0.0000
C(3)	1.244827	0.012055	103.2632	0.0000
R-squared	0.998834	Mean dependent var		2999.436
Adjusted R-squared	0.998822	S.D. dependent var		1459.707
S.E. of regression	50.09460	Akaike info criterion		10.68030
Sum squared resid	504403.2	Schwarz criterion		10.72910
Log likelihood	-1086.391	Hannan-Quinn criter.		10.70004
F-statistics	86081.29	Durbin-Watson stat		0.295995
Prob(F-statistics)	0.000000			

# Nonlinear regression: examples

Examples 7.4 & 7.8 (Greene):

Analysis of a Nonlinear Consumption Function

NLS with starting values equal to the parameters from the OLS estimation ( $c(3)$  equal to 1)

---

Dependent Variable: REALCONS				
Method: Least Squares (Marquard - EViews legacy)				
Sample 1950Q1 2000Q4    Included observations: 204				
Convergence achieved after 80 iterations				
REALCONS=C(1)+C(2)*REALDPI^C(3)				

---

	Coefficient	Std.Error	t-Statistic	Prob.
C(1)	458.7989	22.50149	20.38971	0.0000
C(2)	0.100852	0.010911	9.243447	0.0000
C(3)	1.244827	0.012055	103.2632	0.0000

---

R-squared	0.998834	Mean dependent var		2999.436
Adjusted R-squared	0.998822	S.D. dependent var		1459.707
S.E. of regression	50.09460	Akaike info criterion		10.68030
Sum squared resid	504403.2	Schwarz criterion		10.72910
Log likelihood	-1086.391	Hannan-Quinn criter.		10.70004
F-statistics	86081.28	Durbin-Watson stat		0.295995
Prob(F-statistics)	0.000000			

---

# Quantile regression - LAD

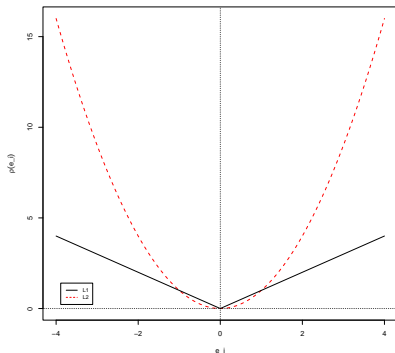
- Quantile regression estimates the relationship between regressors and a specified quantile of dependent variable.
- LAD estimator is the QREG for  $q = \frac{1}{2}$  (median) and the loss function can be described as (compare to OLS objective function):

$$\min_{\hat{\beta}_q} Q_n(\hat{\beta}_q) = \sum_{i=1}^n |y_i - \mathbf{x}'_i \hat{\beta}_q|$$

- LAD estimator predates OLS (itself older than 200 years). Until recently, QREG and LAD have seen little use in econometrics, as OLS is vastly easier to compute.
- Different software packages use a variety of optimization algorithms for QREG/LAD estimation.
- Linear programming can be used for finding QREG estimates (Koenkerr and Bassett (around 1980)).

# Quantile regression - LAD

## OLS vs LAD estimator



Objective (loss) function of the estimators:

$$\min : \sum_{i=1}^N \rho(e_i) = \sum_{i=1}^N \rho(y_i - \mathbf{x}_i' \hat{\beta})$$

For OLS,  $\rho(\cdot)$  is the  $L_2$  norm:  $e_i^2$ .

For LAD,  $\rho(\cdot)$  is the  $L_1$  norm:  $|e_i|$ .

LAD estimation gives much less weight to large deviations.



# Quantile regression - example and motivation

OLS / LAD / QREG coefficient interpretation example:

(1)  $\text{wage}_i = \beta_1 + u_i$

(2)  $\text{wage}_i = \beta_1 + \beta_2 \text{female}_i + u_i$

(3)  $\text{wage}_i = \beta_1 + \beta_2 \text{female}_i + \beta_3 \text{exper}_i + u_i$

The above equations are estimated by OLS / LAD / QREG:

Coefficient	OLS	LAD ( $q = \frac{1}{2}$ )	QREG ( $q = \frac{3}{4}$ )
(1) $\beta_1$	$\hat{\beta}_1 = \bar{y}$ sample mean	$\hat{\beta}_1 = \tilde{y}$ sample median	$\hat{\beta}_1 = Q_3$ sample 3 <sup>rd</sup> quartile
(2) $\beta_1, \beta_1 + \beta_2$	conditional sample mean wage: male / female	cond. sample median wage: male / female	conditional sample $Q_3$ wage: male / female
(3) $\beta_3$	change in expected mean wage for $\Delta \text{exper} = 1$	change in exp. median wage for $\Delta \text{exper} = 1$	change in expected $Q_3$ wage for $\Delta \text{exper} = 1$

# Quantile regression (QREG)

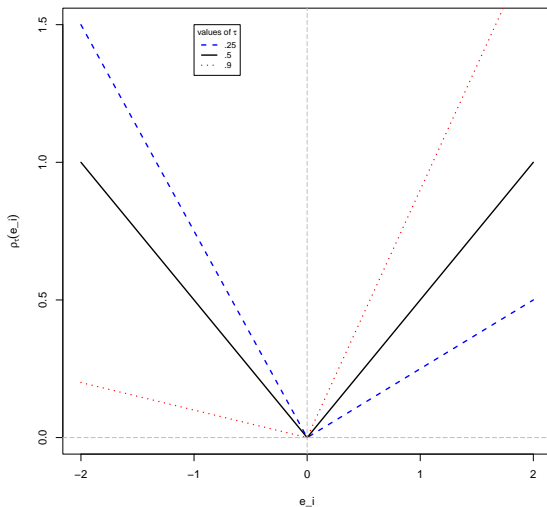
For LRMs, the  $q$ -th quantile QREG estimator  $\beta_q$  minimizes:

$$\min_{\hat{\beta}_q} Q_n(\hat{\beta}_q) = \sum_{i: e_i \geq 0}^n q|y_i - \mathbf{x}'_i \hat{\beta}_q| + \sum_{i: e_i < 0}^n (1 - q)|y_i - \mathbf{x}'_i \hat{\beta}_q|,$$

where  $e_i = (y_i - \mathbf{x}'_i \hat{\beta}_q)$ .

- We use the notation  $\hat{\beta}_q$  to make clear that different choices of  $q$  lead to different  $\hat{\beta}$ .
- Slope of the loss function  $Q_n$  is asymmetrical (around  $e_i = 0$ ).
- The loss function is not differentiable (at  $e_i = 0$ )  
→ gradient methods are not applicable  
(linear programming can be used).

## QREG: different quantiles



For  $q = 0.5$ , positive and negative errors are treated symmetrically.

For  $q \neq 0.5$ , positive and negative errors (of the same magnitude  $e_i$ ) have different weights (penalties), which is reflected in quantile-specific estimates  $\hat{\beta}_q$ .

# Quantile regression (QREG)

- Quantile regression: used to describe relationship between regressors and a specified quantile of dependent variable.
- The (linear) quantile model can be defined as  $Q[y|\mathbf{x}, q] = \mathbf{x}'\beta_q$ , such that  $\text{Prob}[y \leq \mathbf{x}'\beta_q|\mathbf{x}] = q$ ,  $0 < q < 1$  where  $q$  denotes the  $q$ -th quantile of  $y$ .
- One important special case of quantile regression is the least absolute deviations (LAD) estimator, which corresponds to fitting the conditional median of the response variable ( $q = \frac{1}{2}$ ).
- QREG (LAD) estimator can be motivated as a robust alternative to OLS (with respect to outliers).

# Quantile regression example

Example 7.10 (Greene):

Income Elasticity of Credit Cards Expenditure

OLS & LAD & Income elasticity at different deciles

---

Dependent Variable: LOGSPEND				
Method: Least Squares				
Date: 09/15/16 Time 13:53				
Sample (adjusted): 3 13443				
Included observations: 10499 after adjustments				

---

Variable	Coefficient	Std.Error	t-Statistic	Prob.
C	-3.055807	0.239699	-12.74852	0.0000
LOGINC	1.083438	0.032118	33.73296	0.0000
AGE	-0.017364	0.001348	-12.88069	0.0000
ADEPCNT	-0.044610	0.010921	-4.084857	0.0000

---

R-squared	0.100572	Mean dependent var	4.728778
Adjusted R-squared	0.100315	S.D. dependent var	1.404820
S.E. of regression	1.332496	Akaike info criterion	3.412366
Sum squared resid	18634.35	Schwarz criterion	3.415131
Log likelihood	-17909.21	Hannan-Quinn criter.	3.413300
F-statistic	391.1750	Durbin-Watson stat	1.888912
Prob(F-statistic)	0.000000		

---

# Quantile regression example 2

## Example 7.10 (Greene):

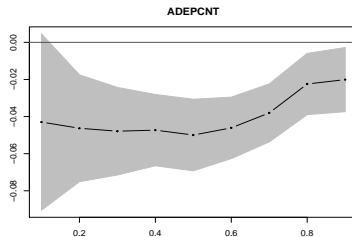
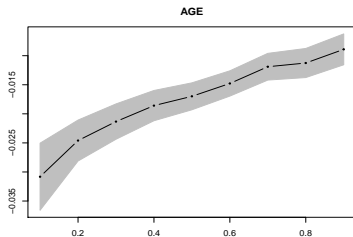
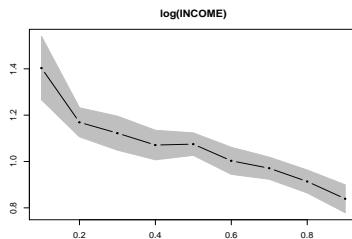
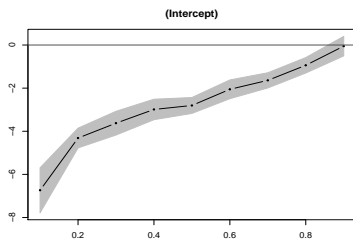
### Income Elasticity of Credit Cards Expenditure (LAD)

Dependent Variable: LOGSPEND    Method: Quantile Regression (Median)				
Sample (adjusted): 3 13443    Included observations: 10499 after adjustments				
Huber Sandwich Standard Errors & Covariance				
Sparsity method: Kernel (Epanechnikov) using residuals				
Bandwidth method: Hall-Sheather, bw=0.04437				
Estimation successfully identifies unique optimal solution				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	-2.803756	0.233534	-12.00577	0.0000
LOGINC	1.074928	0.030923	34.76139	0.0000
AGE	-0.016988	0.001530	-11.10597	0.0000
ADEPCNT	-0.049955	0.011055	-4.518599	0.0000
Pseudo R-squared	0.058243	Mean dependent var		4.728778
Adjusted R-squared	0.057974	S.D. dependent var		1.404820
S.E. of regression	1.346476	Objective		5096.818
Quantile dependent va...	4.941583	Restr. objective		5412.032
Sparsity	2.659971	Quasi-LR statistic		948.0224
Prob(Quasi-LR stat)	0.000000			

# Quantile regression example 2

Example 7.10 (Greene):

Income Elasticity of Credit Cards Expenditure



- Predictions from a CLRM (repetition from BSc courses)
- Predictions: general features,  $k$ FCV, Variance vs. Bias



- CLRM and its estimate:

$$y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots + \beta_K x_K + u$$

$$\hat{y} = \hat{\beta}_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_3 + \cdots + \hat{\beta}_K x_K$$

- Prediction of expected value:

$$\hat{y}_p = E(y|x_1 = 1, x_2 = c_2, \dots, x_K = c_K)$$

$$\hat{y}_p = \hat{\beta}_1 + \hat{\beta}_2 c_2 + \hat{\beta}_3 c_3 + \cdots + \hat{\beta}_K c_K$$

- Rough (underestimated) confidence interval for the expected value prediction: (95%):  $\hat{y}_p \pm 2 \times \text{s.e.}(\hat{y}_p)$ .  
(Rule of thumb)

s.e.( $\hat{y}_p$ ) can be obtained by reparametrization:

- Reparametrized CLRM:

$$y^* = \beta_1^* + \beta_2^*(x_2 - c_2) + \beta_3^*(x_3 - c_3) + \cdots + u$$

- The following holds:

$$\hat{y}_p = \hat{\beta}_1^*$$

$$\text{s.e.}(\hat{y}_p) = \text{s.e.}(\hat{\beta}_1^*), \quad i.e.$$

$$\text{var}(\hat{y}_p) = \text{var}(\hat{\beta}_1^*)$$

- Predicted and actual values of  $y_p$ :

$$\hat{y}_p = \hat{\beta}_1 + \hat{\beta}_2 c_2 + \hat{\beta}_3 c_3 + \cdots + \hat{\beta}_K c_K$$

$$y_p = \beta_1 + \beta_2 c_2 + \beta_3 c_3 + \cdots + \beta_K c_K + u_p$$

- Prediction error

$$\hat{e}_p = y_p - \hat{y}_p = (\beta_1 + \beta_2 c_2 + \beta_3 c_3 + \cdots + \beta_K c_K) + u_p - \hat{y}_p$$

- Prediction error variance

$$\text{var}(\hat{e}_p) = \text{var}(u_p) + \text{var}(\hat{y}_p)$$

because  $\text{var}(\beta_1 + \beta_2 c_2 + \beta_3 c_3 + \cdots + \beta_K c_K) = 0$

- In CLRM, homoscedasticity holds,  $\sigma^2 = \text{var}(u_p)$ :
  - $\text{var}(\hat{e}_p) = \sigma^2 + \text{var}(\hat{y}_p)$
  - We estimate  $\sigma^2$  from the original CLRM as  $(SSR/(n - K))$
  - We get  $\text{var}(\hat{y}_p)$  from the reparametrized LRM
- Standard prediction error:
  - $\text{s.e.}(\hat{e}_p) = \sqrt{\text{var}(\hat{e}_p)}$
- Prediction interval (95%)
  - $\hat{y}_p \pm t_{0.025} \times \text{s.e.}(\hat{e}_p)$

- Prediction with logarithmic dependent variable

$$\log(y) = \beta_1 + \beta_2 x_2 + \cdots + \beta_K x_K + u$$

$$\widehat{\log(y)} = \hat{\beta}_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_K x_K$$

$\hat{y} = e^{\widehat{\log(y)}}$  systematically underestimates  $\hat{y}$  ,

we can use a correction:  $\hat{y} = \hat{\alpha}_0 e^{\widehat{\log(y)}}$

where  $\hat{\alpha}_0 = n^{-1} \sum_{i=1}^n \exp(\hat{u}_i)$

is a consistent (not unbiased) estimator of  $\exp(u)$ .

# Predictions - basics (Matrix form)

Prediction based on estimated model:

$$\hat{y}_p = \mathbf{x}_p' \hat{\boldsymbol{\beta}}$$

Difference between prediction and actual  $y_p$  value:

$$\hat{e}_p = \hat{y}_p - y_p = \mathbf{x}_p' \hat{\boldsymbol{\beta}} - \mathbf{x}_p' \boldsymbol{\beta} - u_p = \mathbf{x}_p' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) - u_p$$

If  $\hat{\boldsymbol{\beta}}$  is unbiased estimator for  $\boldsymbol{\beta}$ ,

$\hat{y}_p$  is an unbiased estimator for  $y_p$  value:

$$E(\hat{e}_p) = E(\hat{y}_p - y_p) = \mathbf{x}_p' E(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + E(-u_p) = 0$$

and the variance of  $\hat{e}_p$  can be expressed as:

$$E(\hat{e}_p^2) = \text{var}(\hat{e}_p) = \mathbf{x}_p' \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_p + \text{var}(u_p)$$

# Predictions - basics (Matrix form)

Variance of  $\hat{e}_p$  (continued):

$$\begin{aligned}\text{var}(\hat{e}_p) &= \mathbf{x}'_p \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{x}_p + \text{var}(u_p) \\ &= \mathbf{x}'_p \left[ \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \right] \mathbf{x}_p + \text{var}(u_p) \\ &\quad \text{substitute } \sigma^2, \text{var}(u_p) \text{ with } \hat{\sigma}^2 \text{ (homoscedasticity)} \\ &= \underbrace{\mathbf{x}'_p \left[ \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} \right] \mathbf{x}_p}_{\hat{\sigma}_p^2} + \hat{\sigma}^2\end{aligned}$$

With growing sample size (asymptotically),

$\text{var}(u_p) = \hat{\sigma}_p^2 + \hat{\sigma}^2$  converges to  $\hat{\sigma}^2$

$\dots \text{plim } \hat{\boldsymbol{\beta}} = \boldsymbol{\beta} \leftrightarrow \text{plim } \hat{\sigma}_p^2 = 0$

(Note: recall consistency of the OLS estimator under A1–A5 conditions & for the CLRM model - i.e. under A1–A6.)

# Predictions - basics (Matrix form)

Variance of  $\hat{e}_p$  (continued):

$$\text{var}(\hat{e}_p) = \mathbf{x}'_p \left[ \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} \right] \mathbf{x}_p + \hat{\sigma}^2$$

after re-arranging, s.e.( $\hat{e}_p$ ) may be written as

$$\text{s.e.}(\hat{e}_p) = \hat{\sigma} \cdot \sqrt{1 + \mathbf{x}'_p (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_p} ,$$

which relates to the individual prediction error.

For mean prediction errors (considering  $\hat{\sigma}_p^2$  only):

$$\text{s.e.}(\tilde{e}_p) = \hat{\sigma} \cdot \sqrt{\mathbf{x}'_p (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_p} .$$

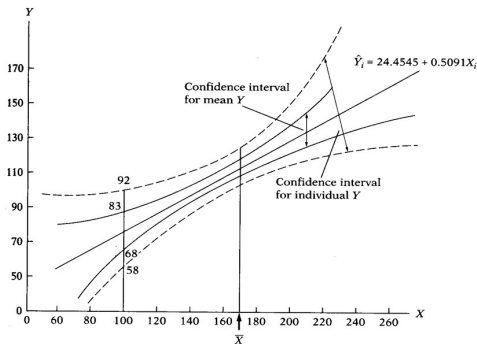


# Predictions - basics (Matrix form)

Prediction intervals: individual vs. mean value predictions:

**Individual prediction:**  $y_p \in \hat{y}_p \pm t_{\alpha/2}^* \times \text{s.e.}(\hat{e}_p)$

**Mean value:**  $y_p \in \hat{y}_p \pm t_{\alpha/2}^* \times \text{s.e.}(\tilde{e}_p)$



# Predictions – general discussion:

- Reliability of predictions:
  - we work with estimated parameters  
(if we generalize from the CLRM paradigm, finite/small sample properties of estimators may be difficult to describe),
  - model parameters can change in time  
(discussed separately in next Block – see Chow tests),
  - predictions include “individual” random errors.
- Impacts of random errors on predictions of individual values are usually much bigger than the impacts of variance in estimated parameters.

# Mean Squared Error of prediction

We can generalize the previous discussion on predictions by considering both biased and unbiased predictors and by allowing for different functional forms and complexity levels in predictive models.

Predictions may be compared/evaluated using:

- $$MSE = E \left[ \left( y_i - \hat{f}(\mathbf{x}_i) \right)^2 \right]$$

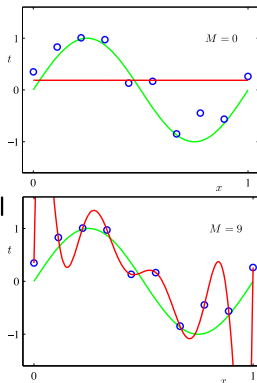
where  $\hat{f}(\mathbf{x}_i)$  is the prediction that  $\hat{f}$  generates for the  $i$ -th regressor set. Here,  $\hat{f}$  represents a general class of predictors (linear, non-linear, non-parametric, etc.) and it may produce either biased or unbiased predictions

# Variance vs. Bias trade-off

Example for a “sine-like” function:  $y = f(x) + u$

## Bias-Variance tradeoff – Intuition

- **Model too simple:** does not fit the data well
  - A *biased* solution
- **Model too complex:** small changes to the data, solution changes a lot
  - A *high-variance* solution



# Train sample & Test sample

Suppose we fit a model  $\hat{f}(\mathbf{x})$  to some training data  $\text{Tr} = \{y_i, \mathbf{x}_i\}_1^n$  and we wish to see how well it performs.

- We could compute  $MSE$  over  $\text{Tr}$ :

$$MSE_{\text{Tr}} = \frac{1}{n} \sum_{i \in \text{Tr}} \left[ y_i - \hat{f}(\mathbf{x}_i) \right]^2$$

When searching for the “best” model by minimizing  $MSE$ , the above statistic would lead to over-fit models.

- Instead, we should (if possible) compute the  $MSE$  using fresh test data  $\text{Te} = \{y_i, \mathbf{x}_i\}_1^m$ :

$$MSE_{\text{Te}} = \frac{1}{m} \sum_{i \in \text{Te}} \left[ y_i - \hat{f}(\mathbf{x}_i) \right]^2$$

## Variance vs. Bias trade-off

Suppose we have a model  $\hat{f}(\mathbf{x})$ , fitted to some training data  $\text{Tr}$  and let  $\{y_0, \mathbf{x}_0\}$  be a test observation drawn from the population. If the true model is  $y_i = f(\mathbf{x}_i) + \varepsilon_i$ , with  $f(\mathbf{x}_i) = E(y_i|\mathbf{x}_i)$ , then the **expected test MSE** can be decomposed into:

$$E(MSE_0) = \text{var}(\hat{f}(\mathbf{x}_0)) + [\text{Bias}(\hat{f}(\mathbf{x}_0))]^2 + \text{var}(\varepsilon_0),$$

where

$$\text{Bias}(\hat{f}(\mathbf{x}_0)) = E[\hat{f}(\mathbf{x}_0)] - f(\mathbf{x}_0),$$

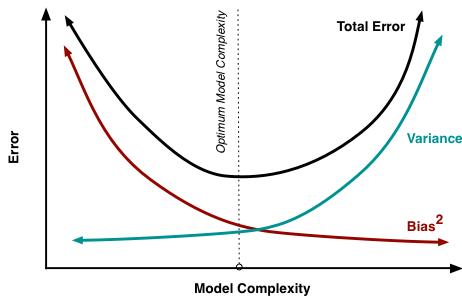
$\varepsilon_0$  is the irreducible error:  $E(MSE_0) \geq \varepsilon_0$ ,

all three RHS elements are non-negative,

The above equation refers to the average test  $MSE$  that we would obtain if we repeatedly estimated  $f(\mathbf{x})$  using a large number of training sets and then tested each  $\hat{f}(\mathbf{x})$  at  $\mathbf{x}_0$ .

# Variance vs. Bias trade-off

$$E(MSE_0) = \text{var}(\hat{f}(\mathbf{x}_0)) + [\text{Bias}(\hat{f}(\mathbf{x}_0))]^2 + \text{var}(\varepsilon_0),$$



This is an illustration,  $\text{var}(\varepsilon_0)$  not shown explicitly.  
(lies at the /asymptotic/ minima of Variance and Bias<sup>2</sup>)

# $k$ -Fold Cross Validation

- Training error ( $MSE_{Tr}$ ) can be calculated easily.
- However,  $MSE_{Tr}$  is not a good approximation for the  $MSE_{Te}$  (out-of sample predictive properties of the model).
- Usually,  $MSE_{Tr}$  dramatically underestimates  $MSE_{Te}$ .

Cross-validation is based on re-sampling (similar to bootstrap).

Repeatedly fit a model of interest to samples formed from the training set & make “test sample” predictions, in order to obtain additional information about predictive properties of the model.

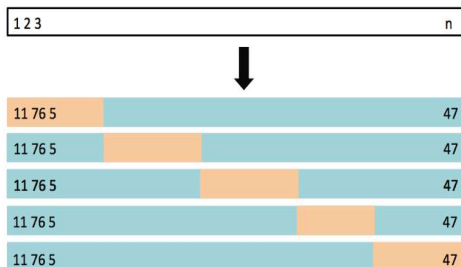


# $k$ -Fold Cross Validation

- In  $k$ -Fold Cross-Validation ( $k$ FCV), the original sample is randomly partitioned into  $k$  roughly equal subsamples (divisibility).
- One of the  $k$  subsamples is retained as the test sample, and the remaining  $(k - 1)$  subsamples are used as training data.
- The cross-validation process is then repeated  $k$  times (the  $k$  folds), with each of the  $k$  subsamples used exactly once as the test sample.
- The  $k$  results from the folds can then be averaged to produce a single estimation.
- $k = 5$  or  $k = 10$  is commonly used.

# $k$ -Fold Cross Validation

$k$ FCV example for CS data &  $k = 5$ :  
(random sampling, no replacement)



In TS, a similar “Walk forward” test procedure may be applied.

# $k$ -Fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{s=1}^k MSE_s,$$

where  $CV_{(k)}$  is the cross-validated estimate of  $MSE$ ,

$k$  is the number of folds used (e.g. 5 or 10),

$$MSE_s = \frac{1}{m_s} \sum_{i \in C_s} (y_i - \hat{y}_i)^2$$

$m_s$  is the number of observations in the  $s$ -th test sample

$C_s$  refers to the  $s$ -th set of test sample observations.

As we evaluate predictions from two or more models,  
we look for the lowest  $CV_{(k)}$ .