

General-to-Specific Time Series Modelling

Felix Pretis

University of Oxford

Michaelmas 2017

Lecture 4: Model Evaluation and Theory of Reduction



Core References for Lecture 4:

- Hendry (2009)* Through the Looking Glass
- Hendry (2017)* Deciding between alternative approaches in macroeconomics
- Doornik (2009) Autometrics
- Pretis, Reade, and Sucarrat (2016) GETS in R
- Hendry (1995) Dynamic Econometrics (as overview)
- Hendry and Doornik (2014) (as overview)



"Even if one wants to test an economic hypothesis as to whether some effect is present, partial inference cannot be conducted alone, unless one is sure about the complete absence of all contaminating influences."

Hendry (2009)

OXFORD MARTIN SCHOOL OXFORD AT THE OXFORD MARTIN SCHOOL

Economies high dimensional, interdependent, heterogeneous, and evolving: comprehensive specification of all events is impossible.

Data generation process (DGP):

economic mechanism plus measurement system.

- Local DGP (LDGP) is for set of variables under analysis.
- Models reflect LDGP not facsimiles:
- designed to satisfy selection criteria.

Aggregation over time, space, commodities, agents, endowments, essential but precludes claim to 'truth'.

Only congruence is on offer in economics: congruent models match LDGP in all measured attributes.

'True' models in class of congruent models.



Congruence is testable: necessary conditions for structure.





 y_t observed when z_t input; $f(\cdot)$ maps inputs to outputs.

 v_t is small, random perturbation.

'Same' outputs repeating experiments at same inputs.

In an econometric model, however:

$y_t = $ [observed]	$g(z_t)$ [explanation]	$+ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \$	(2)
---------------------	---------------------------	---	-----

 y_t decomposed into two components: $g(z_t)$ (bit explained) and ϵ_t (unexplained). Always feasible even if y_t does not depend on $g(z_t)$.

Pretis (Oxford

4: Model Evaluation



In econometrics:

$$\epsilon_{t} = y_{t} - g\left(z_{t}\right)$$

(3)

Thus, models can be **designed** by selection of z_t .

Design criteria must be analyzed: will lead to notion of *congruent* model.

Successive congruent models must explain earlier: concept of **encompassing** – whereby progress achieved.

Will repeatedly use: P(a, b) = P(a | b)P(b)Base order on: time, theory and institutional knowledge.



All random variables in economy over $t=1,\ldots T$: denoted $\{\mathbf{u}_t\}$ with $\mathbf{U}_T^1=(\mathbf{u}_1,\ldots,\mathbf{u}_T)$ Defined on probability space $(\Omega,\mathcal{F},\mathsf{P}).$

DGP is joint data density function $D_U(\cdot)$:

 $\mathsf{D}_{\mathsf{U}}\left(\mathbf{U}_{\mathsf{T}}^{1} \mid \mathbf{U}_{\mathsf{0}}, \boldsymbol{\psi}_{\mathsf{T}}^{1}\right) \text{ with } \boldsymbol{\psi}_{\mathsf{T}}^{1} \in \Psi \subseteq \mathfrak{R}^{\mathsf{k}}, \tag{4}$

- $\psi_{\mathsf{T}}^1 \in \Psi \subseteq \mathfrak{R}^k$ is $k \times 1$ parameter in space Ψ ;
- ψ_{T}^{1} must not depend on \mathcal{F} ;
- U₀ are initial conditions.
- $\mathbf{U}_{\mathrm{T}}^{1}$ unmanageably large: must reduce.

From **DGP** through **LDGP** & **GUM** to selected model arises by sequence of **reductions** organized into twelve stages: see Hendry (1994) and Hendry (2009).

Pretis (Oxford



From DGP to LDGP:

- Aggregation
- 2 Data transformations
- Oata partition
- Marginalizing
- Sequential factorization
- Parameters of interest Approximating the LDGP:
- Lag truncation
- Parameter constancy and invariance
- Functional form (linearity)
 Formulating a specific GUM:
- Mapping to non-integrated data
- Conditional factorization
- Simultaneity

Pretis (Oxford

















1-1 mapping of \mathbf{U}_T^1 to new data set \mathbf{W}_T^1 where $\mathbf{U}_T^1 \leftrightarrow \mathbf{W}_T^1$.

DGP of U_1^1 , and so of W_1^1 , characterized by joint density:

$$\mathsf{D}_{\mathsf{U}}\left(\mathbf{U}_{\mathsf{T}}^{1} \mid \mathbf{U}_{\mathsf{0}}, \psi_{\mathsf{T}}^{1}\right) = \mathsf{D}_{\mathsf{W}}\left(\mathbf{W}_{\mathsf{T}}^{1} \mid \mathbf{W}_{\mathsf{0}}, \phi_{\mathsf{T}}^{1}\right)$$

where $\psi_{\mathsf{T}}^1 \in \Psi$ and $\phi_{\mathsf{T}}^1 \in \Phi$.

Transformation from U to W affects parameter space so Ψ transformed into Φ . But W^{1}_{T} contain aggregates of interest.

Key to all reductions:

• impact on parameters $-\psi_{\rm T}^1$ versus $\phi_{\rm T}^1$.

(5)



Transform joint probabilities to a product:

$$P(a, b) = P(a | b) P(b)$$

= $P(b | a) P(a)$

Only use one ordering (often time); and use repeatedly!

$$P(a, [b, c]) = P(a | [b, c]) P(b, c)$$

= P(a | [b, c]) P(b | c) P(c]

Allowing for parameters, we write:

 $\mathsf{P}(\mathsf{a},\mathsf{b} \mid \boldsymbol{\psi}) = \mathsf{P}(\mathsf{a},\mathsf{b} \mid \boldsymbol{\theta}) = \mathsf{P}(\mathsf{a} \mid \mathsf{b},\boldsymbol{\theta}_1) \mathsf{P}(\mathsf{b} \mid \boldsymbol{\theta}_2)$

where $\boldsymbol{\theta} = \mathbf{f}(\boldsymbol{\psi})$ and $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \in \Theta \subseteq \mathbb{R}^k$.

Crucial issue: does $\Theta = \Theta_1 \times \Theta_2$?







Partition W_{T}^{1} into the two sets:

$$\mathbf{W}_{\mathsf{T}}^1 = \left(\mathbf{X}_{\mathsf{T}}^1: \mathbf{V}_{\mathsf{T}}^1\right)$$

where \mathbf{X}_{T}^{1} is $T \times n$.

Everything desired from the analysis must be learnt from analyzing \mathbf{X}_{T}^{1} alone: \mathbf{V}_{T}^{1} must not be essential to inference.

Hidden condition that parameters of their distributions must be **variation free.**

(6)



$$\begin{array}{ll} \mathsf{D}_{\mathsf{W}}\left(\mathbf{W}_{\mathsf{T}}^{1} \mid \mathbf{W}_{0}, \boldsymbol{\phi}_{\mathsf{T}}^{1}\right) & = & \mathsf{D}_{\mathsf{V}|\mathsf{X}}\left(\mathbf{V}_{\mathsf{T}}^{1} \mid \mathbf{X}_{\mathsf{T}}^{1}, \mathbf{W}_{0}, \boldsymbol{\varphi}_{\mathfrak{a},\mathsf{T}}^{1}\right) \times \\ & & \mathsf{D}_{\mathsf{X}}\left(\mathbf{X}_{\mathsf{T}}^{1} \mid \mathbf{W}_{0}, \boldsymbol{\varphi}_{\mathsf{b},\mathsf{T}}^{1}\right) \end{array}$$

Eliminate $\mathbf{V}_{\mathsf{T}}^{1}$ by **discarding conditional density** $[\mathsf{D}_{\mathsf{V}|\mathsf{X}}\left(\mathbf{V}_{\mathsf{T}}^{1}|\mathbf{X}_{\mathsf{T}}^{1},\mathbf{W}_{0},\varphi_{a,\mathsf{T}}^{1}\right)$ in (7)].

Retain marginal density $D_X \left(\mathbf{X}_T^1 | \mathbf{W}_0, \varphi_{b,T}^1 \right)$.

Key is **no loss of relevant information**: a **cut** is required, so that: $(\varphi_{a,T}^1 : \varphi_{b,T}^1) \in \Phi_a \times \Phi_b$.

Thus, parameters must be variation free.



Innovation process created by sequentially factorizing X1:

$$D_{X} \left(\mathbf{X}_{T}^{1} \mid \mathbf{W}_{0}, \varphi_{b,T}^{1} \right) = D_{x} \left(\mathbf{x}_{T} \mid \mathbf{X}_{T-1}^{1}, \mathbf{W}_{0}, \varphi_{b,T} \right) \times D_{X} \left(\mathbf{X}_{T-1}^{1} \mid \mathbf{W}_{0}, \varphi_{b,T-1}^{1} \right)$$

$$\vdots \qquad (8)$$
$$\mathbf{LDGP} = \prod_{t=1}^{T} D_{x} \left(\mathbf{x}_{t} \mid \mathbf{X}_{t-1}^{1}, \mathbf{W}_{0}, \varphi_{b,t} \right)$$

Creates mean innovation error process: $\epsilon_t = \mathbf{x}_t - E\left[\mathbf{x}_t | \mathbf{X}_{t-1}^1\right]$. Can treat sequential densities 'as if' independent: allows laws of large numbers, central limit theorems etc. as $E\left[\epsilon_t | \mathbf{X}_{t-1}^1\right] = \mathbf{0}$. (implicitly assume $E_t[\cdot] = E[\cdot]$).



 $\begin{array}{l} \mbox{Conditional, sequential distribution of } \{ {\bf x}_t \} \\ \mbox{must not depend on } {\bf V}_{t-1}^1 \mbox{ (Granger non-causality):} \end{array}$

$$\mathsf{D}_{\mathsf{x}}\left(\mathbf{x}_{\mathsf{t}} \mid \mathbf{V}_{\mathsf{t}-1}^{1}, \mathbf{X}_{\mathsf{t}-1}^{1}, \mathbf{W}_{\mathsf{0}}, \phi_{\mathsf{b},\mathsf{t}}\right) = \mathsf{D}_{\mathsf{x}}\left(\mathbf{x}_{\mathsf{t}} \mid \mathbf{X}_{\mathsf{t}-1}^{1}, \mathbf{W}_{\mathsf{0}}, \phi_{\mathsf{b},\mathsf{t}}^{*}\right) \quad (9)$$

If so:

$$\prod_{t=1}^{T} \mathsf{D}_{\mathsf{x}} \left(\mathbf{x}_{t} \mid \mathbf{X}_{t-1}^{1}, \mathbf{W}_{0}, \phi_{b,t}^{*} \right)$$
(10)

from (9) is **LDGP** for $\{\mathbf{x}_t\}$ with $\phi_{b,t}^* = \varphi_{b,t}$ so (9) coincides with last line of (8) above.

LDGP well defined as a reduction of DGP, but may be high-dimensional, complicated, non-linear, & evolving. Knowledge of LDGP is best that can be achieved in space of $\{x_t\}$ -target for model selection.



Denoted $\mu \in \mathcal{M}$. Should be:

- identifiable,
- constant, and
- invariant to relevant class of interventions.

Requires that $\mu = f(\phi_T^1)$ from (5).

Guides choice of:

- data to analyze: \mathbf{X}_{T}^{1} ;
- data transformations: $h(x_t)$;
- form of model: $\mathbf{g}[\mathbf{h}(\mathbf{X}_{\mathsf{T}}^1)] = \boldsymbol{\epsilon}_{\mathsf{t}}$.







Fix extent of history of X_{t-1}^1 in (8) at s earlier periods:

 $\mathsf{D}_{\mathsf{x}}\left(\mathbf{x}_{\mathsf{t}} \mid \mathbf{X}_{\mathsf{t}-1}^{1}, \mathbf{W}_{\mathsf{0}}, \varphi_{\mathsf{b}, \mathsf{t}}\right) = \mathsf{D}_{\mathsf{x}}\left(\mathbf{x}_{\mathsf{t}} \mid \mathbf{X}_{\mathsf{t}-1}^{\mathsf{t}-\mathsf{s}}, \mathbf{W}_{\mathsf{0}}, \varphi_{\mathsf{b}, \mathsf{t}}\right) \tag{11}$

Obvious checks on validity of such a reduction are whether longer lags matter, or error remains an innovation: key criterion is impact on $\{\varphi_{b,t}\}$.



The parameters in question are those that characterize the distribution in (11).

 $\{\varphi_{b,t}\}$ may have different elements for different times:

$$\{\varphi_{b,t}\} = (\varphi_{b,1}, \varphi_{b,2} \cdots \varphi_{b,T-1}, \varphi_{b,T})$$
(12)

Constancy entails that $\{\varphi_{b,t}\}$ depends on a smaller set of parameters that are constant, at least within regimes. Complete parameter constancy is:

$$\varphi_{b,t} = \lambda_0 \in \Lambda_0 \, \forall t.$$
 (13)

E.g.: $y_t = \rho_t z_t + \varepsilon_t$ where $\rho_t = \rho_{t-1} + \eta_t \implies \rho_t = \rho_0 + \sum_{i=1}^t \eta_i$

Most misunderstood, and one of least analyzed, concepts

- not improved by 'random parameters';
- not the same as invariance;
- and non-constancy not entailed by forecast failure.

Pretis (Oxford

4: Model Evaluation



Invariance is under extensions of the information set:

- over time: if extend data series, parameters stay the same;
- across regimes: if change input variable, relation to output variable is unchanged and;
- new sources: if add additional variables to the analysis, parameters are unaltered.

All three are necessary conditions, and are easily testable:

- test for parameter constancy at end of sample;
- Change a regressor and test parameter constancy;
- add potential candidate variables and test they are irrelevant.



Map \mathbf{x}_t into $\mathbf{x}_t^* = \mathbf{h}(\mathbf{x}_t)$.

Denote resulting data by \mathbf{X}^* .

Assume that \mathbf{x}_{t}^{*} makes $D_{x^{*}}(\cdot)$ approximately normal and homoscedastic, $N_{n}[\eta_{t}, \Upsilon]$, with no loss of information:

$$\mathsf{D}_{\mathsf{x}}\left(\mathbf{x}_{\mathsf{t}} \mid \mathbf{X}_{\mathsf{t}-1}^{\mathsf{t}-\mathsf{s}}, \mathbf{W}_{\mathsf{0}}, \boldsymbol{\lambda}_{\mathsf{0}}\right) = \mathsf{D}_{\mathsf{x}^{*}}\left(\mathbf{x}_{\mathsf{t}}^{*} \mid (\mathbf{X}^{*})_{\mathsf{t}-1}^{\mathsf{t}-\mathsf{s}}, \mathbf{W}_{\mathsf{0}}, \boldsymbol{\lambda}\right)$$
(14)

Intimate connection between functional form and parameter constancy.

Drop * notation for simplicity, so \mathbf{x}_t denotes after any relevant transformations.

Linearity



For implementation, must specify precise functional form: e.g., linear model for $D_x(\mathbf{x}_t | \mathbf{X}_{t-1}^{t-s}, \mathbf{W}_0, \boldsymbol{\lambda})$:

$$\mathbf{x}_{t} = \sum_{j=1}^{s} \mathbf{A}_{j} \mathbf{x}_{t-j} + \boldsymbol{\epsilon}_{t}$$
(15)

Express more generally by $\Gamma(L) \mathbf{x}_t = \boldsymbol{\epsilon}_t$ where:

$$\Gamma\left(L\right) = \sum_{j=0}^{s} \Gamma_{j} L^{j};$$

 Γ (L) is a polynomial matrix of order *s* in lag operator L where $\Gamma_0 = I_n$.

Underpins VAR models as the GUM.

Pretis (Oxford)



Most economic time series are non-stationary. Key aspect is unit root (integrated data): denoted I(1). Reduction to I(0) ensures conventional inferences valid but many inferences valid even if I(0) reduction not enforced.

Importantly holds for most diagnostic tests

(not heteroscedasticity however).

Mapping to I(0) also helps interpret outcomes and reduces dimensionality of parameterization. Two possible reductions: **differencing and cointegration**:

- $\Delta \mathbf{x}_t = \mathbf{x}_t \mathbf{x}_{t-1}$ always removes unit root;
- $\beta' \mathbf{x}_t$ sometimes does so as well.

Then both are I(0), and conventional inference applies.



Factorize \mathbf{x}_t into sets of n_1 and n_2 variables, $n_1 + n_2 = n$:

$$\mathbf{x}'_{t} = \left(\mathbf{y}'_{t} : \mathbf{z}'_{t}\right),$$
 (16)

where \mathbf{y}_t endogenous and \mathbf{z}_t non-modelled.

$$\begin{aligned} \mathsf{D}_{\mathsf{x}}\left(\mathbf{x}_{\mathsf{t}} \mid \mathbf{X}_{\mathsf{t}-1}^{\mathsf{t}-\mathsf{s}}, \mathbf{W}_{0}, \boldsymbol{\lambda}\right) &= \mathsf{D}_{\mathsf{y}|\mathsf{z}}\left(\mathbf{y}_{\mathsf{t}} \mid \mathbf{z}_{\mathsf{t}}, \mathbf{X}_{\mathsf{t}-1}^{\mathsf{t}-\mathsf{s}}, \mathbf{W}_{0}, \boldsymbol{\theta}_{\mathfrak{a}}\right) \times \\ \mathsf{D}_{\mathsf{z}}\left(\mathbf{z}_{\mathsf{t}} \mid \mathbf{X}_{\mathsf{t}-1}^{\mathsf{t}-\mathsf{s}}, \mathbf{W}_{0}, \boldsymbol{\theta}_{\mathsf{b}}\right) \end{aligned} \tag{17}$$

 \mathbf{z}_t weakly exogenous for $\boldsymbol{\mu}$ if:

- $\boldsymbol{\mu} = \mathbf{f}(\boldsymbol{\theta}_{a})$ alone; and
- $(\boldsymbol{\theta}_{a}, \boldsymbol{\theta}_{b}) \in \boldsymbol{\Theta}_{a} \times \boldsymbol{\Theta}_{b}.$

Justifies contemporaneous conditioning.

Pretis (Oxford

4: Model Evaluation





Pretis (Oxford)



Consequences of failure of weak exogeneity vary from:

- loss of estimation efficiency,
- through to a loss of parameter constancy.

Experimental setting where Gauss-Markov conditions seem satisfied:

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\beta} + \boldsymbol{\epsilon} \text{ with } \boldsymbol{\epsilon} \sim \mathsf{N}_{\mathsf{T}} \left[\mathbf{0}, \sigma_{\boldsymbol{\epsilon}}^{2} \mathbf{I} \right]$$
 (18)

when $\mathbf{Z}' = (\mathbf{z}_1 \dots \mathbf{z}_T)$ is a $T \times k$ matrix, rank $(\mathbf{Z}) = k$, and $\epsilon' = (\epsilon_1 \dots \epsilon_T)$ with:

$$\mathsf{E}[\mathbf{y} \mid \mathbf{Z}] = \mathbf{Z}\boldsymbol{\beta}$$

hence $E[\mathbf{Z}'\boldsymbol{\epsilon}] = \mathbf{0}$. But:

OLS need not be most efficient unbiased estimator of β .



Explicit weak exogeneity condition required on \mathbb{Z} : must preclude β being learned from marginal distribution. Suppose that marginal equals:

$$\mathbf{z}_{t} = \boldsymbol{\beta} + \boldsymbol{\nu}_{t} \text{ with } \boldsymbol{\nu}_{t} \sim \mathsf{N}_{\mathsf{k}} [\mathbf{0}, \boldsymbol{\Omega}_{\nu}]$$
 (19)

allows sample mean \overline{z} of z to be an unbiased estimator:

$$\mathsf{E}[\bar{\mathbf{z}}] = \boldsymbol{\beta},\tag{20}$$

$$V[\bar{\mathbf{z}}] = \mathsf{T}^{-1} \Omega_{\mathsf{v}},\tag{21}$$

so \overline{z} can dominate $\widehat{\beta}$, possibly dramatically if Ω_{ν} is tiny compared to $\sigma_{\epsilon}^{2}(\mathbf{Z}'\mathbf{Z})^{-1}$.



Cointegrated systems allow testing of one aspect of weak exogeneity:

- equilibrium-correction mechanisms which cross-link equations violate long-run weak exogeneity; also shows weak exogeneity cannot necessarily be obtained merely by choosing 'parameters of interest'.
- The presence of a disequilibrium term in more than one equation is testable.

Structural breaks allow tests for **super exogeneity** and the **Lucas** (1976) critique: see e.g. Engle and Hendry (1993).

When conditional models are constant despite data moments changing, there is evidence of super exogeneity for that model's parameters: automatic test in Hendry and Santos (2010).



Sequentially factorize DGP of n-dimensional $\{x_t\}$:

$$\prod_{t=1}^{T} \mathsf{D}_{\mathsf{x}} \left(\mathbf{x}_{t} \mid \mathbf{X}_{t-1}, \boldsymbol{\theta} \right) = \prod_{t=1}^{T} \mathsf{D}_{\mathsf{y}|\mathsf{z}} \left(\mathbf{y}_{t} \mid \mathbf{z}_{t}, \mathbf{X}_{t-1}, \boldsymbol{\phi}_{1} \right) \mathsf{D}_{\mathsf{z}} \left(\mathbf{z}_{t} \mid \mathbf{X}_{t-1}, \boldsymbol{\phi}_{2} \right)$$
(22)
where $\mathbf{x}_{t} = \left(\mathbf{y}_{t}' : \mathbf{z}_{t}' \right)$ and $\boldsymbol{\phi} = \left(\boldsymbol{\phi}_{1}' : \boldsymbol{\phi}_{2}' \right)' = \mathbf{f} \left(\boldsymbol{\theta} \right) \in \mathbb{R}^{k}$.

If parameters of **y** and **z** processes variation free – \mathbf{z}_t weakly exogenous for parameters of interest $\boldsymbol{\psi} = \mathbf{h}(\phi_1)$: does not rule out that ϕ_1 may change if ϕ_2 is changed.

Super exogeneity adds parameter invariance in conditional:

$$\frac{\partial \phi_1}{\partial \phi_2'} = \mathbf{0}.$$



If $D_{x}(\cdot)$ is multivariate normal:

$$\left(\begin{array}{c} y_t \\ z_t \end{array} \right) \sim IN_n \left[\left(\begin{array}{c} \mu_{1,t} \\ \mu_{2,t} \end{array} \right), \left(\begin{array}{c} \sigma_{11,t} & \sigma'_{12,t} \\ \sigma_{12,t} & \sigma_{22,t} \end{array} \right) \right] -$$

 $\mu_{1,t}$ and $\mu_{2,t}$ are functions of X_{t-1} . Suppose economic theory suggests that $\mu_{1,t} = \beta \mu_{2,t}$ β is primary parameter of interest.

Cond. Model:
$$\mathsf{E}[\mathsf{y}_t|\mathsf{z}_t] = \mu_{1,t} + \sigma_{12,t}\sigma_{22,t}^{-1}(\mathbf{z_t} - \mu_{2,t})$$

Cond. Var.:
$$V[y_t|z_t] = \sigma_{11,t} - \sigma'_{12,t}\sigma_{22,t}^{-1}\sigma_{12,t} = \omega_t^2$$

e cond. model as:

$$\mathsf{E}[\mathsf{y}_t|\mathsf{z}_t] = \beta \mu_{2,t} + \sigma_{12,t} \sigma_{22,t}^{-1} \mathbf{z}_t - \sigma_{12,t} \sigma_{22,t}^{-1} \mu_{2,t}$$

when is z_t super exogenous for parameters (β , ω_t^2)?

Pretis (Oxford

re-writ



Further re-write cond. model:

$$\begin{split} \mathsf{E}[\mathsf{y}_t|\mathsf{z}_t] &= \ \beta \mu_{2,t} + \sigma_{12,t} \sigma_{22,t}^{-1} \mathsf{z}_t - \sigma_{12,t} \sigma_{22,t}^{-1} \mu_{2,t} \\ &= \ (\beta - \sigma_{12,t} \sigma_{22,t}^{-1}) \mu_{2,t} + \sigma_{12,t} \sigma_{22,t}^{-1} \mathsf{z}_t \\ &= \ \gamma_{1,t} + \gamma_{2,t} \mathsf{z}_t \end{split}$$

Econometric model:

$$y_t = \beta_0 + \beta z_t + \varepsilon_t$$

For z_t to be weakly exogenous:

• $\beta = \sigma_{12} \sigma_{22}^{-1} \forall t$ then $\mu_{2,t}$ does not enter conditional model.

For z_t to be super exogenous additionally need:

•
$$\gamma_{2,t} = \gamma_2 \forall t$$

• $\omega_t^2 = \omega^2 \forall t$

• $(\gamma_1, \gamma_2, \omega^2)$ invariant to changes in marginal $(\mu_{2,t}, \sigma_{22,t})$.



Key implication: $\gamma_{1,t}$ depends on $\mu_{2,t}$ when $\beta \neq \sigma_{12}\sigma_{22}^{-1} = \gamma_{2,t}$. Then changes in marginal reflected in conditional. requires $\beta = \gamma_2$, which is testable if $\mu_{2,t}$ shifts. If γ_2 unaffected by shifts in $\mu_{2,t}$, then \mathbf{z}_t super exogenous.

Note: Automatic test using indicators from marginal tested for significance in conditional.



Three attributes:

- 'uniqueness';
- 'corresponds to desired entity';
- 'satisfies the assumed interpretation'.

Consider regression of quantity on price and want to 'identify as demand equation':

- unique function of data second moments;
- but need not correspond to underlying demand behaviour;
- and may be incorrectly interpreted-really supply schedule.

Uniqueness is often the sense intended in econometrics.



Simultaneous representation is a reduction.

$$B \mathbf{y}_t = C \mathbf{z}_t + \boldsymbol{\epsilon}_t \ \ \text{where} \ \ \boldsymbol{\epsilon}_t \sim \mathsf{IN}_{n_1}\left[0, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}}\right] \tag{24}$$
 so actual DGP is:

$$\mathbf{y}_{t} = \mathbf{\Pi} \mathbf{z}_{t} + \mathbf{w}_{t}$$
 where $\mathbf{w}_{t} \sim \mathsf{IN}_{n_{1}}[\mathbf{0}, \mathbf{\Omega}_{w}]$ (25)

with:

$$\mathbf{B}\mathbf{\Pi} - \mathbf{C} = \mathbf{0} \tag{26}$$

(26) may entail restrictions on Π , so is a testable reduction.



$$\begin{split} \text{GUM:} \ \mathbf{A}\left(\mathbf{L}\right)\mathbf{h}\left(\mathbf{y}\right)_{t} = \mathbf{B}\left(\mathbf{L}\right)\mathbf{g}\left(\mathbf{z}\right)_{t} + \epsilon_{t} \end{split}$$
 where $\epsilon_{t \ \widetilde{app}} \ \mathsf{N}_{n_{1}}\left[\mathbf{0}, \boldsymbol{\Sigma}_{\varepsilon}\right].$

Notice the many design steps in reaching (27).

 ϵ_t is a **derived** and not an autonomous process:

$$\varepsilon_{t} = \mathbf{A} (\mathbf{L}) \mathbf{h} (\mathbf{y})_{t} - \mathbf{B} (\mathbf{L}) \mathbf{g} (\mathbf{z})_{t}$$
(28)

Designed by choices made in reduction.

Crucial insight: model changes as reductions are altered.

Parameters of model are **functions** of ϕ_T^1 in (5): $(\mathbf{A}(\cdot) \mathbf{B}(\cdot)) = \mathbf{f}(\phi_T^1)$ alter as ϕ_T^1 changes.

Test empirical model against LDGP, as represented by data properties.

(27)



Measures of no information loss if correct

- Aggregation: no loss of information on marginalizing wrt disaggregates if retain sufficient statistics for μ;
- **Transformations:** no associated reduction, but introduce μ , and need for parameters to be **invariant** and **identifiable**;
- Data partition: which variables to include/omit fundamental to success of empirical modelling;
- Marginalizing wrt v_t without loss, if X¹_T sufficient for μ marginalizing wrt V¹_{t-1} without loss, if Granger non-causality for x_t and a cut;
- Sequential factorization: no loss if ε_t an innovation relative to X¹_{t-1}



- LDGP: reduction of DGP for relevant variables, nested within it; properties explained by reduction;
- Lag truncation: no loss if et remains an innovation;
- Parameter constancy and invariance: constancy across interventions on marginal process;
- Functional form: no reduction when two densities in (14) are equal
- Integrated data:

reduce to I(0) by **cointegration** and **differencing** Conventional inference and more parsimonious;

- Conditional factorization: eliminate marginal process
 No loss of information if z_t weakly exogenous for μ;
- Simultaneity: parsimoniously capture joint dependence.



 Knowledge of DGP entails knowledge of all reductions thereof
 If one model entails knowledge of others, then that model is said to encompass them

Unlike symptomatology approach testing for:

autocorrelation, heteroscedasticity, omitted variables, non-constancy etc. then 'correcting' them.

Invalid approach as no unique alternative to any null; **successive outcomes can contradict**;

no obvious termination pointstopping after first 'non-rejection' is disastrous, and no account of selection process.



Commence from 'good' approximation to LDGP, which embeds available economic theory, and institutional knowledge, checked for congruence by mis-specification tests.

Mimic reduction theory in practical research, to minimize losses due to reductions imposed

Based on notions of **congruence** and **encompassing**:

- former close to 'well specified',
- latter entails explaining the results of other models.

We will address how to select model of LDGP from initial general unrestricted specification



Partition own data \mathbf{X}_{T}^{1} into three information sets:

- past data;
- Present data;
- future data;

$$\mathbf{X}_{\mathsf{T}}^{1} = \left(\mathbf{X}_{t-1}^{1} : \mathbf{x}_{t} : \mathbf{X}_{\mathsf{T}}^{t+1}\right) \tag{29}$$

- **Theory information**: source of μ , and creative stimulus;
- Measurement information: price index theory, identities, data accuracy; and:
- Rival models (encompass congruent models)

Leads to six main model evaluation criteria.



- homoscedastic innovation ϵ_t ;
- 2 weakly exogenous \mathbf{z}_t for $\boldsymbol{\mu}$;
- (a) constant, invariant μ ;
- theory consistent, identifiable structures;
- Idata-admissible formulations on accurate observations;
- encompass rival models.

Exhaustive nulls to test; but many alternatives.

Models which satisfy [1] & [2] are **well specified** on available information.

Models which satisfy [1] & [2] & [3] are (empirically) congruent.

Admissible, theory-consistent, encompassing, congruent model satisfies all six criteria.



- Define a starting model: general unrestricted model (GUM)
 - Designed to be congruent (diagnostic testing) and relevant,
 - Tests of reductions with approximately correct distribution,
 - Reduction can maintain congruence (or lack thereof),
 - Reduction up to a predefined significance level (backtesting w.r.t. GUM: acceptable information loss).



- Define a starting model: general unrestricted model (GUM)
 - Designed to be congruent (diagnostic testing) and relevant,
 - Tests of reductions with approximately correct distribution,
 - Reduction can maintain congruence (or lack thereof),
 - Reduction up to a predefined significance level (backtesting w.r.t. GUM: acceptable information loss).

Model selection is an iterative search procedure, need to follow several paths:

- multiple path search, or
- tree search.



- t-tests (single variable removal).
- F-tests (tests of variables removed from the GUM, encompassing aka backtesting).
- F-tests (pruning to faster search).
- diagnostic tests
 - ARCH (Engle 1982)
 - Serial correlation (Godfrey 1978, Harvey 1981)
 - Heteroscedasticity (White 1980)
 - Normality (Jarque and Bera 1980; Doornik and Hansen 1994, 2008)
 - Chow (Chow 1960 in-sample stability test)*
- information criterion (tiebreaker)
- stability tests (out of sample, optionally)



- Model selection is an iterative search procedure
 - manual search can follow a few paths: slow and tedious,
 - computer automated search can follow all paths, Well, not all. There are 2^k models, so need a strategy. k = 100 at 10^9 /sec: $10^6 \times$ age of universe.



- Model selection is an iterative search procedure
 - manual search can follow a few paths: slow and tedious,
 - computer automated search can follow all paths, Well, not all. There are 2^k models, so need a strategy. k = 100 at 10^9 /sec: $10^6 \times$ age of universe.
- General-to-specific model selection (Gets, 'Hendry' or 'LSE' methodology) largely driven by David Hendry (DHSY, PcGive, Alchemy, Dynamic Econometrics, ...) Lively debate.
- Automated Gets initiated by Hoover and Perez (1999), Hendry and Krolzig (2005) (PcGets: 2nd generation, theoretical properties, bias correction).
 Study model selection through simulation – improves debate.
- Autometrics & Getsm improve on PcGets, extended beyond standard regression models.



- Hoover-Perez (1999):
 - General unrestricted model
 - 2 Multiple path search
 - Encompassing test
 - Diagnostic testing
 - 5 Tiebreaker
- Hendry and Krolzig (1999), PcGets (2001):
 - Add presearch
 - 2 Extend multiple-path search
 - Add iteration
 - No out-of-sample testing (Lunch and Vital-Ahuja, 1998)
 - Ohange treatment for Invalid GUM
- Autometrics (2009), Getsm (2017):
 - Reduce role of presearch
 - 2 Change search path algorithm: tree search
 - Extend scope: separation of model and algorithm
 - Increase efficiency



Autometrics & Getsm implement underlying principle of general-to-specific model selection (*'Hendry methodology'*).

Autometrics & Getsm

- likelihood-based
- searches the whole model space:
 - tree search ensures that no model is estimated twice
 - irrelevant paths can be cut-off efficiently
 - F-tests are used to speed-up search
- implements backtracking on diagnostics: only test from terminal candidates, then backtrack if necessary
- backtesting w.r.t. GUM 0 (the initial GUM after presearches) removes need for encompassing of candidate models
- relevant terminal candidates remembered in iterated search
- $\bullet\,$ implements block search for $N\geqslant T$





Search follows branches till no insignificant variables; tests for congruence and parsimonious encompassing; backtracks if either fails, till first non-rejection found.



Path search gives impression of 'repeated testing'. Confused with selecting from 2^{N} possible **models** (here $2^{1000} = 10^{301}$, an impossible task). We are selecting **variables**, not models, & only N variables.

But selection matters, as only retain 'significant' outcomes. Sampling variation also entails retain irrelevant, or miss relevant, by chance near selection margin.

Conditional on selecting, estimates biased away from origin: but can bias correct as know c_{α} .

Small efficiency cost under null for examining many candidate regressors, even N >> T.

Almost as good as commencing from LDGP at same c_{α} .



DGP -

Model -

$$\begin{array}{rcl} \mathfrak{X}_{\text{fixed}} &=& \{1\}\\ \mathfrak{X}_{\text{free}} &=& \{y_{t-1}, x_{1t}, x_{2t}, x_{3t}, x_{4t}\} \end{array}$$

Four selection methods





gauge: fraction of irrelevant variables (x_{3t}, x_{4t}) in the final model potency: fraction of relevant variables $(y_{t-1}, x_{1t}, x_{2t})$ in the final model.

Pretis (Oxford



Diagnostic checking: subject every estimated model to a battery of diagnostic tests (normality, residual correlation, residual ARCH, in-sample Chow, out-of-sample Chow,...). If a test fails the reduction is rejected and the next model in line is considered.

Under null of congruent GUM, the figure below compares 'gauges' for Autometrics with diagnostic checking **on** vs. **off**:

$$y_t = \sum\nolimits_{i=1}^N \beta_i z_{i,t} + \varepsilon_t \ \text{ for } \ \varepsilon_t \sim \text{IN}[0,\sigma_\varepsilon^2] \eqno(30)$$

 $T = 100, n = 1, \dots, 10 = N; \beta_k = 0 \text{ for } k > n; R^2 = 0.9.$

Gauge is close to α if diagnostic tests **not** checked.

Gauge is larger than α with diagnostics **on**, when checking to ensure a congruent reduction.

Difference seems due to retaining insignificant irrelevant variables which proxy chance departures from null of mis-specification tests.

Pretis (Oxford

4: Model Evaluation





Pretis (Oxford)







Motivated by Indicator Saturation Block-Partitioning.

More Variables Than Observations: Block Partitioning

```
\begin{array}{l} \mbox{Split variables into blocks } A, \ B, \ C \\ A \cup B \ \mbox{select} \rightarrow G_1 \\ A \cup C \ \mbox{select} \rightarrow G_2 \\ B \cup C \ \mbox{select} \rightarrow G_3 \\ \Rightarrow \ G_1 \cup G_2 \cup G_3 \ \mbox{select} \rightarrow \mbox{Final Model} \end{array}
```



If also have relevant variables to be retained, and N > T, orthogonalize them with respect to the rest.

As N > T, divide in more sub-blocks, setting $\alpha = 1/N$.

Basic model retains desired sub-set of **n** variables at every stage, and only selects over putative irrelevant variables at stringent significance level:

under the null, has no impact on estimated coefficients of relevant variables, or their distributions.

Thus, almost costless to check even large numbers of candidate variables:

huge benefits if initial specification incorrect but enlarged GUM nests LDGP.



Hoover and Perez (1999) experiments: $HP7y_{7,t} = 0.75y_{7,t-1} + 1.33x_{11,t} - 0.9975x_{11,t-1} + 6.44u_t \quad R^2 = 0.58$ $HP8y_{8,t} = 0.75y_{8,t-1} - 0.046x_{3,t} + 0.0345x_{3,t-1} + 0.073\lambda u_t R^2 = 0.93$ where $u_t \sim IN[0, 1]$; $x_{i,t-j}$ are US macro data 1) The GUM has 3 DGP variables plus 37 irrelevant. 2) Then consider 141 irrelevant, larger than T = 139.



N < T: T = 139, 3 relevant and 37 irrelevant variables

	Hoover-Perez		step-wise		Autometrics		
	HP7	HP8	HP7	HP8	HP7	HP8	
	1% nominal size						
Gauge %	3.0*	0.9*	0.9	3.1	1.6	1.6	
Potency %	94.0	99.9	100.0	53.3	99.2	100.0	
DGP found %	24.6	78.0	71.6	22.0	68.3	68.8	

* Only counting significant terms (but tiebreaker is best-fitting model)

N > T: T = 139, 3 relevant and 141 irrelevant variables

	step-wise		Autometrics		
	HP7	HP8	HP7	HP8	
	0.1% nominal size				
Gauge %	0.1	0.7	0.3	0.1	
Potency %	99.7	40.3	97.4	100.0	
DGP found %	87.4	9.0	82.9	90.2	

Oxford)

4: Model Evaluation



Have now established that there is little loss from using the path-search algorithms

- Gauge is close to selected α for both.
- Potency is near theory value for a 1-off test.
- Goodness-of-fit is not directly used to select models & no attempt is made to 'prove' that a given set of variables matters, and 'repeated testing' is not a concern, but choice of c_{α} affects \mathbb{R}^2 and n through retention by $|t_{(n)}| \ge c_{\alpha}$.
- Likelihood estimation in general is feasible (Doornik (2009)).



Commence with general formulation – general unrestricted model:

$$y_t = \beta' \mathbf{z}_t + \gamma' \mathbf{w}_t + \sum_{j=1}^T \delta_{\text{IIS},j} \mathbf{1}_{\{j=t\}} + \sum_{j=1}^{T-1} \delta_{\text{SIS},j} \mathbf{1}_{\{j \leqslant t\}} + \nu_t \quad t = 1, \dots, T$$

- Embed theory z_t
- Expand model w_t (almost costless if theory correct)
- Indicators δ_t (almost costless under null)
- Ensuring valid conditioning exogeneity
 - Testable (Weak E., Strong E., Super E.)
- Theory Motivation: Reduction from DGP to LDGP to GUM to Specific



Doornik, J. A. (2009). Autometrics. pp. 88–121. Oxford: Oxford University Press.

Engle, R. F. and D. F. Hendry (1993). Testing super exogeneity and invariance in regression models. Journal of Econometrics 56, 119–139. Reprinted in Ericsson, N. R. and Irons, J. S. (eds.) Testing Exogeneity, Oxford: Oxford University Press, 1994.

Hendry, D. F. (1995). Dynamic Econometrics. Oxford: Oxford University Press.

Hendry, D. F. (2009).
The methodology of empirical econometric modeling: Applied econometrics through the looking-glass.
In T. C. Mills and K. D. Patterson (Eds.), *Palgrave Handbook of Econometrics*, pp. 3–67. Basingstoke: Palgrave MacMillan.

Hendry, D. F. (2017). Deciding between alternative approaches in macroeconomics. International Journal of Forecasting. Forthcoming.

Hendry, D. F. and J. A. Doornik (2014). Empirical Model Discovery and Theory Evaluation. Cambridge, Mass.: MIT Press.

Hendry, D. F. and H.-M. Krolzig (2005). The properties of automatic Gets modelling. Economic Journal 115, C32–C61.

Pretis (Oxford)



Hendry, D. F. and C. Santos (2010).

An automatic test of super exogeneity.

In M. W. Watson, T. Bollerslev, and J. Russell (Eds.), Volatility and Time Series Econometrics, pp. 164–193. Oxford: Oxford University Press.

Hoover, K. D. and S. J. Perez (1999).

Data mining reconsidered: Encompassing and the general-to-specific approach to specification search.

Econometrics Journal 2, 167-191.

Lucas, R. E. (1976). Econometric policy evaluation: A critique.

In K. Brunner and A. Meltzer (Eds.), The Phillips Curve and Labor Markets, Volume 1 of Carnegie-Rochester Conferences on Public Policy, pp. 19–46. Amsterdam: North-Holland Publishing Company.

Pretis, F., J. Reade, and G. Sucarrat (2016).

General-to-specific (gets) modelling and indicator saturation with the r package gets.

Oxford Department of Economics Discussion Paper 794.