

Cross-sectional and spatial dependence in panels with **R**

Giovanni Millo¹

¹Research Dept., Assicurazioni Generali S.p.A.

Innsbruck
April 25th 2012



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels



An applied example: Munnell 1990, again...

Munnell (1990), *Public capital productivity*:

Does public capital (roads, water facilities, public buildings and structures) help growth?

48 US states, annual data 1970-1986

Production function:

$$\log(gdp) = \alpha + \beta_1 \log(pcap) + \beta_2 \log(pc) + \beta_3 \log(emp) + \beta_4 unemp$$



The Munnell model by OLS

Simple OLS estimation, taken at face value, tells you:

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	1.643302	0.057411	28.6237	< 2.2e-16	***
log(pcap)	0.155007	0.017101	9.0641	< 2.2e-16	***
log(pc)	0.309190	0.010240	30.1930	< 2.2e-16	***
log(emp)	0.593935	0.013705	43.3362	< 2.2e-16	***
unemp	-0.006733	0.001412	-4.7683	1.858e-06	***

At 95% significance, 1% increase in public capital raises GDP by 0.12-0.19%



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → inconsistency
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, inconsistency
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → inconsistency
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, inconsistency
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → inconsistency
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, inconsistency
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → inconsistency
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, inconsistency
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → inconsistency
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, inconsistency
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → **inconsistency**
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, **inconsistency**
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → **inconsistency**
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, **inconsistency**
- nonstationarity → spurious regressions



Can we rely on this result?

...it depends on the underlying hypotheses.

- Coefficients are unbiased if regressors are exogenous
- Inference is valid if errors are (asymptotically) normal...
- ...and uncorrelated.

Panel data are prone to a number of specification issues:

- unobserved individual heterogeneity
 - correlated with regressors → **inconsistency**
 - uncorrelated with the regressors → inefficiency
- serial correlation → inefficiency
- cross-sectional/spatial correlation → inefficiency, **inconsistency**
- nonstationarity → spurious regressions

Applied example 2 - Enter Baltagi!

Standard panel problem: unobserved heterogeneity $u_{it} = \mu_i + \epsilon_{it}$

Adding random effects (here, by ML):

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
$\text{phi} = \sigma_{\mu}^2 / \sigma_{\epsilon}^2$	5.0005	1.0972	4.5577	5.172e-06	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.14386583	0.13440520	15.9508	< 2.2e-16	***
log(pcap)	0.00314439	0.02348562	0.1339	0.8935	
log(pc)	0.30981115	0.01991177	15.5592	< 2.2e-16	***
log(emp)	0.73133720	0.02502053	29.2295	< 2.2e-16	***
unemp	-0.00613818	0.00090629	-6.7729	1.262e-11	***



Applied example 2 - Enter Baltagi!

Standard panel problem: unobserved heterogeneity $u_{it} = \mu_i + \epsilon_{it}$
 Adding random effects (here, by ML):

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
$\text{phi} = \sigma_{\mu}^2 / \sigma_{\epsilon}^2$	5.0005	1.0972	4.5577	5.172e-06	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.14386583	0.13440520	15.9508	< 2.2e-16	***
log(pcap)	0.00314439	0.02348562	0.1339	0.8935	
log(pc)	0.30981115	0.01991177	15.5592	< 2.2e-16	***
log(emp)	0.73133720	0.02502053	29.2295	< 2.2e-16	***
unemp	-0.00613818	0.00090629	-6.7729	1.262e-11	***



Applied example 2 - Enter Baltagi!

Standard panel problem: unobserved heterogeneity $u_{it} = \mu_i + \epsilon_{it}$
 Adding random effects (here, by ML):

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
$\text{phi} = \sigma_{\mu}^2 / \sigma_{\epsilon}^2$	5.0005	1.0972	4.5577	5.172e-06	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.14386583	0.13440520	15.9508	< 2.2e-16	***
log(pcap)	0.00314439	0.02348562	0.1339	0.8935	
log(pc)	0.30981115	0.01991177	15.5592	< 2.2e-16	***
log(emp)	0.73133720	0.02502053	29.2295	< 2.2e-16	***
unemp	-0.00613818	0.00090629	-6.7729	1.262e-11	***



2 sides to the talk:

Robustness features against XS correlation in the `plm` package

- XS-dependence without any explicit spatial characteristic (e.g., due to the presence of common factors)
- OLS/FE/RE estimates are still consistent but for valid inference we need robust covariance matrices
- Pesaran's augmentation approach allows for (nonstationary) common factors and valid unit root tests

Spatial models characterizing XS dependence in a parametric way (`splm` package)

- explicitly taking distance into account
- distance matrix is exogenous and time-invariant (although it needn't be *geographic* distance)
- the estimation framework is ML or GM



2 sides to the talk:

Robustness features against XS correlation in the `p1m` package

- XS-dependence without any explicit spatial characteristic (e.g., due to the presence of common factors)
- OLS/FE/RE estimates are still consistent but for valid inference we need robust covariance matrices
- Pesaran's augmentation approach allows for (nonstationary) common factors and valid unit root tests

Spatial models characterizing XS dependence in a parametric way (`sp1m` package)

- explicitly taking distance into account
- distance matrix is exogenous and time-invariant (although it needn't be *geographic* distance)
- the estimation framework is ML or GM



2 sides to the talk:

Robustness features against XS correlation in the `p1m` package

- XS-dependence without any explicit spatial characteristic (e.g., due to the presence of common factors)
- OLS/FE/RE estimates are still consistent but for valid inference we need robust covariance matrices
- Pesaran's augmentation approach allows for (nonstationary) common factors and valid unit root tests

Spatial models characterizing XS dependence in a parametric way (`sp1m` package)

- explicitly taking distance into account
- distance matrix is exogenous and time-invariant (although it needn't be *geographic* distance)
- the estimation framework is ML or GM



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

Why R

R is Open Source: you are

- free to use it
- free to modify and redistribute it
- free to look at the code (and see what it's doing)

R is scalable (a matrix language at heart)

R is extendable (*"useRs become developeRs"*)

R has the best graphics in the stats world

R has a wealth of flexible and transparent optimizers

Check it out yourself: www.r-project.org



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels



The plm package for panel data econometrics

The plm package for panel data Econometrics (see JSS 27/2):

- provides infrastructure for panel data::
 - special data structure: the `pdata.frame`
 - special `pseries` class for individual series
 - convenient lag and difference methods for `pseries`
 - special formula interface for common data manipulations (*demeaning*)
 - support for extended Formulas (IV)
- includes a number of estimators and tests

In particular, it implements the general framework of robust restriction testing (see package `sandwich`, Zeileis, JStatSoft 2006) based upon correspondence between conceptual and software tools in

$$W = (R\beta - r)' [R'vcov(\beta)R]^{-1} (R\beta - r)$$



The plm package for panel data econometrics

The plm package for panel data Econometrics (see JSS 27/2):

- provides infrastructure for panel data::
 - special data structure: the `pdata.frame`
 - special `pseries` class for individual series
 - convenient lag and difference methods for `pseries`
 - special formula interface for common data manipulations (*demeaning*)
 - support for extended Formulas (IV)
- includes a number of estimators and tests

In particular, it implements the general framework of robust restriction testing (see package `sandwich`, Zeileis, JStatSoft 2006) based upon correspondence between conceptual and software tools in

$$W = (R\beta - r)' [R'vcov(\beta)R]^{-1} (R\beta - r)$$



The plm package for panel data econometrics

The plm package for panel data Econometrics (see JSS 27/2):

- provides infrastructure for panel data::
 - special data structure: the `pdata.frame`
 - special `pseries` class for individual series
 - convenient lag and difference methods for `pseries`
 - special formula interface for common data manipulations (*demeaning*)
 - support for extended Formulas (IV)
- includes a number of estimators and tests

In particular, it implements the general framework of robust restriction testing (see package `sandwich`, Zeileis, JStatSoft 2006) based upon correspondence between conceptual and software tools in

$$W = (R\beta - r)'[R'vcov(\beta)R]^{-1}(R\beta - r)$$



Robustness features for panel models

Robust covariance estimators are based on the White framework, a.k.a. the *sandwich* estimator. The `plm` version of the robust covariance estimator `vcovHC()` is based on White's formula and (partial) demeaning

$$vcov(\beta) = (X'X)^{-1} \sum_h X_h E_h X_h' (X'X)^{-1}$$

We need a `vcov` estimator robust vs. XS correlation. 3 possibilities:

- White cross-section: $E_t = e_t e_t'$ is robust w.r.t. arbitrary heteroskedasticity and XS-correlation; depends on T-asymptotics
- Beck & Katz unconditional XS-correlation (a.k.a. PCSE): $E = \frac{\sum_i \epsilon_t \epsilon_t'}{N}$
- or the Driscoll and Kraay (RES 1998) estimator, mixing kernel-weighted sums of $\epsilon_t \epsilon_{t-l}'$, robust vs. time-space correlation decreasing in time ...



Robust diagnostic testing under XSD

... and the trick of robust diagnostic testing is done! Just supply the relevant `vcov` to `coeftest{lmtest}` or `linear.hypothesis{car}`

```
> coeftest(zz, vcov=function(x) vcovHC(x, cluster="time"))
```

	Estimate	Std. Error	t-value	p-value	
log(pcap)	-0.0261497	0.0454291	-0.5756	0.56504	
log(pc)	0.2920069	0.0479729	6.0869	1.821e-09	***
log(emp)	0.7681595	0.0627143	12.2486	< 2.2e-16	***
unemp	-0.0052977	0.0015224	-3.4799	0.00053	***

```
> coeftest(zz, vcov=vcovSCC)
```

	Estimate	Std. Error	t-value	p-value	
log(pcap)	-0.0261497	0.0575413	-0.4545	0.6496338	
log(pc)	0.2920069	0.0588387	4.9628	8.57e-07	***
log(emp)	0.7681595	0.0828411	9.2727	< 2.2e-16	***
unemp	-0.0052977	0.0014912	-3.5528	0.0004046	***



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels



Testing for XS dependence

The CD test 'family' (Breusch-Pagan 1980, Pesaran 2004) is based on transformations of the product-moment correlation coefficient of a model's residuals, defined as

$$\hat{\rho}_{ij} = \frac{\sum_{t=1}^T \hat{u}_{it} \hat{u}_{jt}}{(\sum_{t=1}^T \hat{u}_{it}^2)^{1/2} (\sum_{t=1}^T \hat{u}_{jt}^2)^{1/2}}$$

and comes in different flavours appropriate in N-, NT- and T- asymptotic settings:

$$CD = \sqrt{\frac{2T}{N(N-1)}} \left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N \hat{\rho}_{ij} \right)$$

$$LM = \sum_{i=1}^{N-1} \sum_{j=i+1}^N T_{ij} \hat{\rho}_{ij}^2$$

$$SCLM = \sqrt{\frac{1}{N(N-1)}} \left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N \sqrt{T_{ij}} \hat{\rho}_{ij}^2 \right)$$

Friedman's (1928) rank test and Frees' (1995) test substitute Spearman's rank coefficient for ρ

A CD test of the Munnell model

```
> pcdtest(zz)
```

Pesaran CD test for cross-sectional dependence in panels

data: formula $z = 30.3685$, p-value $< 2.2e-16$

alternative hypothesis: cross-sectional dependence

```
> pcdtest(log(gsp) ~ log(pcap) + log(pc) + log(emp) +
unemp, data=Produc)
```

Pesaran CD test for cross-sectional dependence in panels

data: formula $z = 40.1977$, p-value $< 2.2e-16$

alternative hypothesis: cross-sectional dependence



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 **Digression: testing CSD vs. CWD**
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

Cross-sectional and spatial dependence in panels

Panel data (*longitudinal* data) have a double dimension:

- one is usually (but not always) time
- the other may be individuals, families, groups, firms, countries...

If the cross-sectional dimension has some form of ordering, then we can use *spatial* methods. Space can be

- physical
- economic
- social
- ...

(geographic space has a nice feature: it is exogenous)

Tobler's Law

(Spatial) panel data may exhibit dependence

- over time...
- ...or through the cross-sectional dimension.

Is cross-sectional dependence, if any, related to proximity in (this particular) space?

Tobler's First Law of Geography :-) "Everything is related to everything else, but near things are more related than distant things (Tobler 1970)"



Weak and strong cross-sectional dependence

In recent work by Pesaran and Tosetti dependence in spatial processes is characterized as being *distance-decaying* or not in an asymptotic fashion, introducing the concepts of cross-sectional strong dependence (CSD) and cross-sectional weak dependence (CWD).

Factor models and, respectively, spatial models are typical cases of the first and second type. To fix ideas, we identify the factor model

$$y_{it} = X_{it}\beta + \gamma_i\mu_t + \epsilon_{it}$$

with CSD, and the spatial (autoregressive) error model

$$y_{it} = X_{it}\beta + u_{it}; \quad u_{it} = \lambda(I_T \otimes W)u_{it} + \epsilon_{it}$$

with CWD.



CSD and-or CWD? (And why bother?)

Estimators consistent in the presence of both have been devised (Pesaran's CCE for example); yet one might be interested in the informative content of the dependence structure, not only in the β s

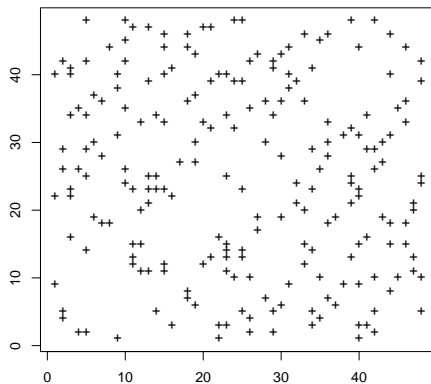
- The economic meaning of CSD and CWD is different:
 - idiosyncratic reaction to unobserved common factor
 - ...vs. spatial diffusion process
- Tests for spatial dependence have power against factor-type dependence and the reverse (*how much power?*)
- If a CSD and a CWD process coexist, the first prevails (asymptotically)

So how do we tell?



Introducing georeferentiation: the local CD tests (1)

Restricting the test to *neighbouring* observations: meet the W matrix!
(here: binary proximity matrix, 48 continental US states)



The local CD tests (2)

The CD(p) test is *CD* restricted to neighbouring observations

$$CD = \sqrt{\frac{T}{\sum_{i=1}^{N-1} \sum_{j=i+1}^N w(p)_{ij}}} \left(\sum_{i=1}^{N-1} \sum_{j=i+1}^N [w(p)]_{ij} \hat{\rho}_{ij} \right)$$

where $[w(p)]_{ij}$ is the (i, j) -th element of the p -th order proximity matrix, so that if h, k are not neighbours, $[w(p)]_{hk} = 0$ and $\hat{\rho}_{hk}$ gets "killed"; W is employed here as a binary selector: any matrix coercible to boolean will do

`pcdtest(..., w=W)` will compute the local test. Else if `w=NULL` the global one.

Only CD(p) is documented, but in principle any of the above tests (LM, SCLM, Friedman, Frees) can be restricted.

A *local* CD test of the Munnell model

```
> pcdtest(zz, w=usaww)
```

Pesaran CD test for cross-sectional dependence in panels

data: formula $z = 17.5521$, p-value $< 2.2e-16$

alternative hypothesis: cross-sectional dependence

```
> pcdtest(log(gsp) ~ log(pcap) + log(pc) + log(emp) +
unemp, data=Produc, w=usaww)
```

Pesaran CD test for cross-sectional dependence in panels

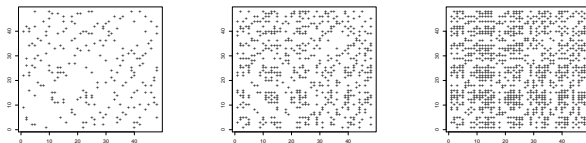
data: formula $z = 17.2397$, p-value $< 2.2e-16$

alternative hypothesis: cross-sectional dependence



Recursive $CD(p)$

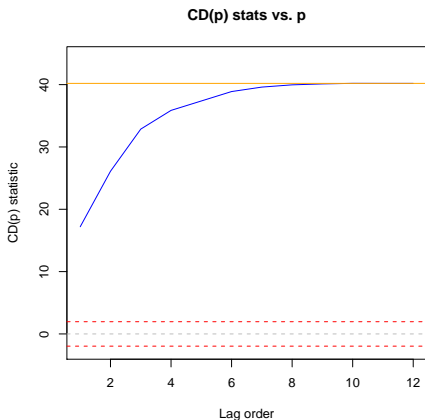
The $CD(p)$ on orders greater than 1 uses higher-order proximity matrices (neighbours of neighbours and so on). Here we can see proximity matrices of the USA states "filling up" as p increases from 1 to 3 (saturation would be at 11).



For a given set of residuals the $CD(p)$ test must converge to CD as p increases, although there is no guarantee it will do so monotonically.

Recursive $CD(p)$ plot for the Munnell model

The CD test, seen as a descriptive statistic, can provide an informal assessment of the degree of 'localness' of the dependence: let the neighbourhood order p grow until $CD(p) \rightarrow CD$



The power problem of CD plots

This might seem a good way to assess whether the dependence is *distance-decaying*. Unfortunately the sample size on which calculation is based is clearly depending on p , and so will test power be.

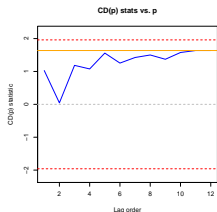
Simulations in Pesaran and Tosetti show that the $CD(p)$ test statistic has lower values than CD if dependence is indeed CSD. The reason for this is that under global dependence the $CD(p)$ turns out to be based on fewer observations (the neighbours) w.r.t. the CD , which exploits the full sample.

Thus the CD has maximum power against CSD while under CWD, on the converse, the $CD(p)$ uses fewer observations, but is likely to use those which are more strongly correlated (the neighbours), possibly offsetting the former effect. Which one will prevail will depend on the degree of spatial correlation.

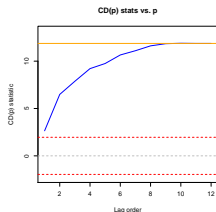


Different patterns for the CD(p) test

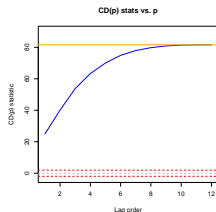
no CWD, no CSD



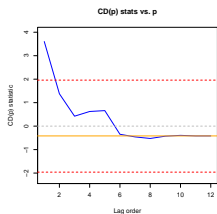
no CWD, weak CSD



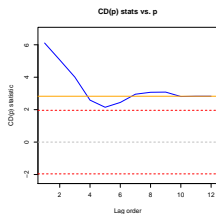
no CWD, strong CSD



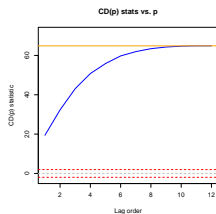
weak CWD, no CSD



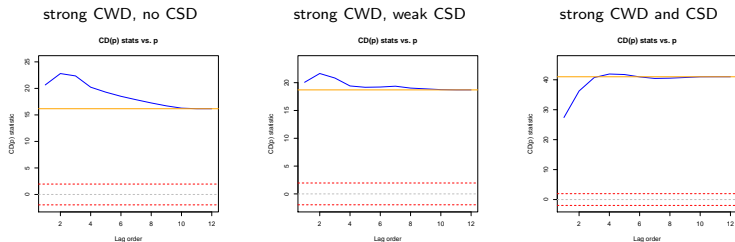
weak CWD and CSD



weak CWD, strong CSD



Different patterns for the $CD(p)$ test, cont.d



- If we assume that there is either one or the other type of dependence, while a decreasing sequence of $CD(p)$ tests is unambiguously indicative of a CWD process, an increasing one might well be related to CWD with *local* correlation approaching unity as well as to CSD.
- If both CWD and CSD coexist, then the overall process is (asymptotically) CSD but the situation as regards the test in finite samples should be less clear-cut.



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models**
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

Spatial models

Standard tool: the *spatial lag operator* $L(x) = Wx$ where W is a *proximity matrix*. In the simplest case, W is a binary (standardized) *neighbourhood indicator* and $L(x)$ is the average of neighbours.

Spatial econometric models have either a spatially lagged dependent variable or error \rightarrow **feedback** \rightarrow **nonlinearity** \rightarrow **endogeneity?**

The two standard specifications:

- Spatial Lag (SAR): $y = \psi W_1 y + X\beta + \epsilon$
- Spatial Error (SEM): $y = X\beta + u; u = \lambda W_2 u + \epsilon$

The general model (Anselin 1988):

$$y = \psi W_1 y + X\beta + u; u = \lambda W_2 u + \epsilon; E[\epsilon\epsilon'] = \Omega$$

Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels



The general estimation framework (Anselin 1988)

Define $A = I - \psi W_1$ and $B = I - \lambda W_2$ (a.k.a. *spatial filters*), then the general log-likelihood is

$$\log L = -\frac{N}{2} \ln \pi - \frac{1}{2} \ln |\Omega| + \ln |A| + \ln |B| - \frac{1}{2} e' e$$

The likelihood is thus a function of β, ψ, λ and parameters in Ω . The overall errors' covariance can be scaled as $B' \Omega B = \sigma_e^2 \Sigma$. This likelihood can be concentrated w.r.t. β and σ_e^2 substituting $e = [\hat{\sigma}_e^2 \Sigma]^{-\frac{1}{2}} (Ay - X\hat{\beta})$

$$\log L = -\frac{N}{2} \ln \pi - \frac{N}{2} \hat{\sigma}_e^2 - \frac{1}{2} \ln |\Sigma| + \ln |A| - \frac{1}{2 \hat{\sigma}_e^2} (Ay - X\hat{\beta})' \Sigma^{-1} (Ay - X\hat{\beta})$$

(ML for correlated errors as in Magnus (1978) plus spatial filter on y)



The general estimation framework - contd.

A closed-form GLS solution for β and σ_e^2 is available for any given set of spatial parameters ψ, λ and scaled covariance matrix Σ

$$\begin{aligned}\hat{\beta} &= (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}Ay \\ \hat{\sigma}_e^2 &= \frac{(Ay - X\hat{\beta})'\Sigma^{-1}(Ay - X\hat{\beta})}{N}\end{aligned}\quad (1)$$

so that a two-step procedure is possible which alternates optimization of the concentrated likelihood and GLS estimation.

(Oberhofer and Kmenta 1974, Magnus 1978)

Operationalizing the general estimation method

The general estimation method can be made operational for specific Σ s parameterized as $\Sigma(\theta)$ by plugging in the relevant Σ , Σ^{-1} and $|\Sigma|$ into the log-likelihood and then optimizing by a two-step procedure, alternating:

$$GLS : \beta = (X'[\Sigma(\hat{\theta})^{-1}]X)^{-1}X'[\Sigma(\hat{\theta})^{-1}]Ay \rightarrow \hat{\beta}$$

and

$$ML : \max_{\theta} ll(\theta|\hat{\beta}) \rightarrow \hat{\theta}$$

until convergence

(yes it does converge: see Oberhofer and Kmenta (*Econometrica* 1974))

The computational problem

The computational problem:

- $\Sigma = \Sigma(\theta)$
- $A = A(\psi)$

so all inverses and determinants are to be recomputed at every optimization loop

Anselin (ibid.) gives efficient procedures for estimating the "simple" cross-sectional SAR and SEM specifications (SAR has a closed-form solution, SEM doesn't): see package `spdep` by Roger Bivand for very fast R versions.

There are few software implementations for more general models (notably, Matlab routines by Elhorst (IRSR 2003) for FE/RE SAR/SEM panels).



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

R for spatial econometrics

A lot of infrastructure is available for spatial econometrics and statistics:
see the relevant *Task View* on cran.r-project.org

CRAN Task View: Analysis of Spatial Data

- Classes for spatial data
- Handling spatial data
- Reading and writing spatial data
- Point pattern analysis
- Geostatistics
- Disease mapping and areal data analysis
- Spatial regression
- Ecological analysis

See also *R-News* 1/2, 1/3, 2/2



R for spatial *panel* econometrics

Estimators and tests for spatial panel models in package **splm** by Gianfranco Piras and me

- ML and GM estimation of spatial panels with:
 - random or fixed effects
 - spatial lag (SAR), spatial error (SEM) or both (SARAR)
- (mis)specification testing
 - spatial Hausman test for FE vs. RE
 - Baltagi et al. conditional tests for effects
 - restriction testing

See

<https://r-forge.r-project.org/projects/splm>

now on **CRAN** as version 1.0-0 and in *JStatSoft* **47(1)**



Applied example 3 - Munnell vindicated?

Add random effects *and SEM*

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	7.536915	1.721237	4.3788	1.193e-05	***
rho	0.558002	0.033401	16.7061	< 2.2e-16	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3366185	0.1383335	16.8912	< 2.2e-16	***
log(pcap)	0.0440546	0.0219709	2.0051	0.0449490	*
log(pc)	0.2437035	0.0200338	12.1646	< 2.2e-16	***
log(emp)	0.7447614	0.0241424	30.8487	< 2.2e-16	***
unemp	-0.0037777	0.0010623	-3.5561	0.0003764	***



Applied example 3 - Munnell vindicated?

Add random effects *and SEM*

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	7.536915	1.721237	4.3788	1.193e-05	***
rho	0.558002	0.033401	16.7061	< 2.2e-16	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3366185	0.1383335	16.8912	< 2.2e-16	***
log(pcap)	0.0440546	0.0219709	2.0051	0.0449490	*
log(pc)	0.2437035	0.0200338	12.1646	< 2.2e-16	***
log(emp)	0.7447614	0.0241424	30.8487	< 2.2e-16	***
unemp	-0.0037777	0.0010623	-3.5561	0.0003764	***



Applied example 3 - Munnell vindicated?

Add random effects *and SEM*

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	7.536915	1.721237	4.3788	1.193e-05	***
rho	0.558002	0.033401	16.7061	< 2.2e-16	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3366185	0.1383335	16.8912	< 2.2e-16	***
log(pcap)	0.0440546	0.0219709	2.0051	0.0449490	*
log(pc)	0.2437035	0.0200338	12.1646	< 2.2e-16	***
log(emp)	0.7447614	0.0241424	30.8487	< 2.2e-16	***
unemp	-0.0037777	0.0010623	-3.5561	0.0003764	***



Applied example 3 - ... or not!? Add SAR

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	21.3175	8.3133	2.5643	0.01034	*

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
lambda	0.161615	0.029108	5.5522	2.821e-08	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3736012	0.1394745	17.0182	< 2.2e-16	***
log(pcap)	0.01294505	0.02493997	0.5190	0.6037	
log(pc)	0.22555376	0.02163422	10.4258	< 2.2e-16	***
log(emp)	0.67081075	0.02642113	25.3892	< 2.2e-16	***
unemp	-0.00579716	0.00089175	-6.5009	7.984e-11	***



Applied example 3 - ... or not!? Add SAR

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	21.3175	8.3133	2.5643	0.01034	*

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
lambda	0.161615	0.029108	5.5522	2.821e-08	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3736012	0.1394745	17.0182	< 2.2e-16	***
log(pcap)	0.01294505	0.02493997	0.5190	0.6037	
log(pc)	0.22555376	0.02163422	10.4258	< 2.2e-16	***
log(emp)	0.67081075	0.02642113	25.3892	< 2.2e-16	***
unemp	-0.00579716	0.00089175	-6.5009	7.984e-11	***



Applied example 3 - ... or not!? Add SAR

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
<code>phi</code>	21.3175	8.3133	2.5643	0.01034	*

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
<code>lambda</code>	0.161615	0.029108	5.5522	2.821e-08	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3736012	0.1394745	17.0182	< 2.2e-16	***
<code>log(pcap)</code>	0.01294505	0.02493997	0.5190	0.6037	
<code>log(pc)</code>	0.22555376	0.02163422	10.4258	< 2.2e-16	***
<code>log(emp)</code>	0.67081075	0.02642113	25.3892	< 2.2e-16	***
<code>unemp</code>	-0.00579716	0.00089175	-6.5009	7.984e-11	***



Applied example 3 - SAREM (SARAR) model

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
<code>phi</code>	7.530808	1.749436	4.3047	1.672e-05	***
<code>rho</code>	0.536835	0.035164	15.2666	< 2.2e-16	***

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value
<code>lambda</code>	0.0018174	0.0091512	0.1986	0.8426

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3736012	0.1394745	17.0182	< 2.2e-16	***
<code>log(pcap)</code>	0.0425013	0.0222146	1.9132	0.055721	.
<code>log(pc)</code>	0.2415077	0.0202971	11.8987	< 2.2e-16	***
<code>log(emp)</code>	0.7419074	0.0244212	30.3797	< 2.2e-16	***
<code>unemp</code>	-0.0034560	0.0010605	-3.2589	0.001119	**



Applied example 3 - SAREM (SARAR) model

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
<code>phi</code>	7.530808	1.749436	4.3047	1.672e-05	***
<code>rho</code>	0.536835	0.035164	15.2666	< 2.2e-16	***

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value
<code>lambda</code>	0.0018174	0.0091512	0.1986	0.8426

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3736012	0.1394745	17.0182	< 2.2e-16	***
<code>log(pcap)</code>	0.0425013	0.0222146	1.9132	0.055721	.
<code>log(pc)</code>	0.2415077	0.0202971	11.8987	< 2.2e-16	***
<code>log(emp)</code>	0.7419074	0.0244212	30.3797	< 2.2e-16	***
<code>unemp</code>	-0.0034560	0.0010605	-3.2589	0.001119	**



Applied example 3 - SAREM (SARAR) model

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	7.530808	1.749436	4.3047	1.672e-05	***
rho	0.536835	0.035164	15.2666	< 2.2e-16	***

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value
lambda	0.0018174	0.0091512	0.1986	0.8426

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	2.3736012	0.1394745	17.0182	< 2.2e-16	***
log(pcap)	0.0425013	0.0222146	1.9132	0.055721	.
log(pc)	0.2415077	0.0202971	11.8987	< 2.2e-16	***
log(emp)	0.7419074	0.0244212	30.3797	< 2.2e-16	***
unemp	-0.0034560	0.0010605	-3.2589	0.001119	**



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

A more general (panel) model with serial correlation

Let us consider a panel model within a more general setting, allowing for a spatially lagged response and the following features of the composite error term (i.e., parameters describing Σ):

- random effects ($\phi = \sigma_{\mu}^2 / \sigma_{\epsilon}^2$)
- spatial correlation in the idiosyncratic error term (λ)
- serial correlation in the idiosyncratic error term (ρ)

$$y = \psi(I_T \otimes W_1)y + X\beta + u$$

$$u = (\iota_T \otimes \mu) + \epsilon$$

$$\epsilon = \lambda(I_T \otimes W_2)\epsilon + \nu$$

$$\nu_t = \rho\nu_{t-1} + e_t$$

Available models

Suitable zero restrictions to the spatial lag and covariance parameters give rise to the following possibilities:

$par \neq 0$	$\mu\psi\rho$	$\mu\rho$	$\mu\psi$	$\psi\rho$	ρ	ψ	μ	(none)
λ	SAREMSRRE	SAREMRE	SARSRRE	SAREMSR	SAREM	SARSR	SARRE	SAR
(none)	SEMSRRE	SEMRE	SRRE	SEMSR	SEM	SR	RE	OLS

where SARRE, SEMRE are the 'usual' random effects spatial panels and SAR, SEM the standard spatial models (here, pooling with $W = I_T \otimes w$)

- The models without serial correlation are already available in the **splm** package: see function `spml()` (used to be `sprem1()`, now hidden by this general wrapper for ML estimation)
- The models in red are new in the literature; those among them including serial correlation are not visible in the user space yet

Baltagi et al.'s LM testing framework

Most applications concentrate on the error model. In this setting, Baltagi et al. (2007) derive conditional LM tests for

- $\lambda | \rho, \mu$ (needs SRRE estimates of \hat{u})
- $\rho | \lambda, \mu$ (needs SEMRE estimates of \hat{u})
- $\mu | \lambda, \rho$ (needs SEMSR estimates of \hat{u})

So a viable and computationally parsimonious strategy for the error model can well be to test in the three directions by means of conditional LM tests and see whether one can estimate a simpler model than the general one. (also forthcoming in **splm**)

An asymptotically equivalent test, although much heavier on the machine, is the Wald test implicit in the diagnostics of the general model. The lag specification can be tested for only the second way. (the covariance for spatial lag and error covariance parameters is based on the numerical estimate of the Hessian).



Performance and stability issues

Optimizing over this relatively populous parameter space is a computationally heavy task which can nevertheless be accomplished by modern computers, even the lowest-end ones for moderately sized samples. Efficient optimizing routines are available in **R**:

- `nlminb()` (PORT)
- `optim()`

and will soon be chosen from at user level. Further issues to consider:

- determinants (are not *that* important up to moderately big N)
- error covariances (are *very* important: efficient use of block structure is crucial for speed; specialized methods for *bds* matrices)
- starting values are important for stability and reliability (not for speed)
- parameter scaling in optimization is crucial to reduce the problem of negative variance estimates for non-significant parameters



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

Applied example 4 - The whole picture?

... or at least, a more complete one! Add **serial correlation** as well.

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	9.5931838	11.5634706	0.8296	0.4068	
psi	0.9883539	0.0042231	234.0342	<2e-16	***
rho	0.6327643	0.0297443	21.2734	<2e-16	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	3.04678280	0.19577778	15.5625	< 2.2e-16	***
log(pcap)	0.04012277	0.03319499	1.2087	0.226778	
log(pc)	0.07020168	0.02141021	3.2789	0.001042	**
log(emp)	0.91163172	0.03016620	30.2203	< 2.2e-16	***
unemp	-0.00245537	0.00074879	-3.2791	0.001041	**



Applied example 4 - The whole picture?

... or at least, a more complete one! Add **serial correlation** as well.

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	9.5931838	11.5634706	0.8296	0.4068	
psi	0.9883539	0.0042231	234.0342	<2e-16	***
rho	0.6327643	0.0297443	21.2734	<2e-16	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	3.04678280	0.19577778	15.5625	< 2.2e-16	***
log(pcap)	0.04012277	0.03319499	1.2087	0.226778	
log(pc)	0.07020168	0.02141021	3.2789	0.001042	**
log(emp)	0.91163172	0.03016620	30.2203	< 2.2e-16	***
unemp	-0.00245537	0.00074879	-3.2791	0.001041	**



Applied example 4 - The whole picture?

... or at least, a more complete one! Add **serial correlation** as well.

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	9.5931838	11.5634706	0.8296	0.4068	
psi	0.9883539	0.0042231	234.0342	<2e-16	***
rho	0.6327643	0.0297443	21.2734	<2e-16	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	3.04678280	0.19577778	15.5625	< 2.2e-16	***
log(pcap)	0.04012277	0.03319499	1.2087	0.226778	
log(pc)	0.07020168	0.02141021	3.2789	0.001042	**
log(emp)	0.91163172	0.03016620	30.2203	< 2.2e-16	***
unemp	-0.00245537	0.00074879	-3.2791	0.001041	**



Applied example 4 - Lag or error? or both?

Fit the full model in order to discriminate between lag and error:

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	9.4499833	26.3098562	0.3592	0.7195	
psi	0.9884254	0.0074188	133.2334	<2e-16	***
rho	0.6263348	0.0424575	14.7520	<2e-16	***

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value
lambda	0.0083648	0.0340687	0.2455	0.806

Coefficients: (*much the same, omitted*)

```
spreml ← splm::spreml
system.time(flmod ← spreml(fm, Produc, w=usaww, errors='semsrre', lag=T))
```

user	system	elapsed
82.097	9.232	91.538



Applied example 4 - Lag or error? or both?

Fit the full model in order to discriminate between lag and error:

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	9.4499833	26.3098562	0.3592	0.7195	
psi	0.9884254	0.0074188	133.2334	<2e-16	***
rho	0.6263348	0.0424575	14.7520	<2e-16	***

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value
lambda	0.0083648	0.0340687	0.2455	0.806

Coefficients: (*much the same, omitted*)

```
spreml ← splm::spreml
system.time(flmod ← spreml(fm, Produc, w=usaww, errors='semsrre', lag=T))
```

user	system	elapsed
82.097	9.232	91.538



Applied example 4 - Lag or error? or both?

Fit the full model in order to discriminate between lag and error:

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
phi	9.4499833	26.3098562	0.3592	0.7195	
psi	0.9884254	0.0074188	133.2334	<2e-16	***
rho	0.6263348	0.0424575	14.7520	<2e-16	***

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value
lambda	0.0083648	0.0340687	0.2455	0.806

Coefficients: (*much the same, omitted*)

```
spreml ← splm::spreml
system.time(flmod ← spreml(fm, Produc, w=usaww, errors='semsrre', lag=T))
```

user	system	elapsed
82.097	9.232	91.538



Pre-transforming the data (here: into first differences)

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
psi	0.253728	0.036909	6.8744	6.224e-12	***
rho	-0.114104	0.073107	-1.5608	0.1186	

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
lambda	0.156042	0.040724	3.8317	0.0001273	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
log(pcap)	0.03925713	0.04388949	0.8945	0.3711	
log(pc)	0.00755859	0.02149750	0.3516	0.7251	
log(emp)	0.90669288	0.03350830	27.0588	< 2.2e-16	***
unemp	-0.00439363	0.00075112	-5.8494	4.933e-09	***



Pre-transforming the data (here: into first differences)

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
psi	0.253728	0.036909	6.8744	6.224e-12	***
rho	-0.114104	0.073107	-1.5608	0.1186	

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
lambda	0.156042	0.040724	3.8317	0.0001273	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
log(pcap)	0.03925713	0.04388949	0.8945	0.3711	
log(pc)	0.00755859	0.02149750	0.3516	0.7251	
log(emp)	0.90669288	0.03350830	27.0588	< 2.2e-16	***
unemp	-0.00439363	0.00075112	-5.8494	4.933e-09	***



Pre-transforming the data (here: into first differences)

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
psi	0.253728	0.036909	6.8744	6.224e-12	***
rho	-0.114104	0.073107	-1.5608	0.1186	

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
lambda	0.156042	0.040724	3.8317	0.0001273	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
log(pcap)	0.03925713	0.04388949	0.8945	0.3711	
log(pc)	0.00755859	0.02149750	0.3516	0.7251	
log(emp)	0.90669288	0.03350830	27.0588	< 2.2e-16	***
unemp	-0.00439363	0.00075112	-5.8494	4.933e-09	***



Pre-transforming the data (here: into first differences)

Error variance parameters:

	Estimate	Std. Error	t-value	p-value	
psi	0.253728	0.036909	6.8744	6.224e-12	***
rho	-0.114104	0.073107	-1.5608	0.1186	

Spatial autoregressive coefficient:

	Estimate	Std. Error	t-value	p-value	
lambda	0.156042	0.040724	3.8317	0.0001273	***

Coefficients:

	Estimate	Std. Error	t-value	p-value	
log(pcap)	0.03925713	0.04388949	0.8945	0.3711	
log(pc)	0.00755859	0.02149750	0.3516	0.7251	
log(emp)	0.90669288	0.03350830	27.0588	< 2.2e-16	***
unemp	-0.00439363	0.00075112	-5.8494	4.933e-09	***



Or ...maybe, a *spatial between* model?

Truer to Munnell's spirit (long-term perspective), a model on time averages (*between*) can be estimated controlling for possible spatial correlation of residuals by extracting the between data with **plm** and estimating the XS model with `errorsarlm()` in **spdep**.

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	1.3621274	0.2218737	6.1392	8.294e-10	***
log(pcap)	0.1485725	0.0645968	2.3000	0.02145	*
log(pc)	0.3759930	0.0422845	8.8920	< 2.2e-16	***
log(emp)	0.5469648	0.0558546	9.7927	< 2.2e-16	***
unemp	-0.0113285	0.0095923	-1.1810	0.23760	

Lambda: 0.55044, LR test value: 7.5521, p-value: 0.005994

But: *what is the implied spatial diffusion process for the errors?*



Or ...maybe, a *spatial between* model?

Truer to Munnell's spirit (long-term perspective), a model on time averages (*between*) can be estimated controlling for possible spatial correlation of residuals by extracting the between data with **plm** and estimating the XS model with `errorsarlm()` in **spdep**.

Coefficients:

	Estimate	Std. Error	t-value	p-value	
(Intercept)	1.3621274	0.2218737	6.1392	8.294e-10	***
log(pcap)	0.1485725	0.0645968	2.3000	0.02145	*
log(pc)	0.3759930	0.0422845	8.8920	< 2.2e-16	***
log(emp)	0.5469648	0.0558546	9.7927	< 2.2e-16	***
unemp	-0.0113285	0.0095923	-1.1810	0.23760	

Lambda: 0.55044, LR test value: 7.5521, p-value: 0.005994

But: *what is the implied spatial diffusion process for the errors?*



Solution: I might have chosen a better example

There are a number of issues with Munnell's model on levels

- variables are probably $I(1)$, residuals also "close to" having a unit root: persistence in time leads to spurious regression a la Granger and Newbold (1974)
- bad identification of random effect vs. serial correlation coefficient is evident (Calzolari and Magazzini 2010)

that can be solved by first-differencing the data. Yet, as Munnell herself (JEL 1992) observes, **first differencing destroys the economic significance of the model.**

So far, the moral of the story is:

- specification analysis of a general model highlighted the limits of our empirical investigation
- panel time series methods might give the correct perspective



Solution: I might have chosen a better example

There are a number of issues with Munnell's model on levels

- variables are probably $I(1)$, residuals also "close to" having a unit root: persistence in time leads to spurious regression a la Granger and Newbold (1974)
- bad identification of random effect vs. serial correlation coefficient is evident (Calzolari and Magazzini 2010)

that can be solved by first-differencing the data. Yet, as Munnell herself (JEL 1992) observes, **first differencing destroys the economic significance of the model.**

So far, the moral of the story is:

- specification analysis of a general model highlighted the limits of our empirical investigation
- panel time series methods might give the correct perspective



Solution: I might have chosen a better example

There are a number of issues with Munnell's model on levels

- variables are probably $I(1)$, residuals also "close to" having a unit root: persistence in time leads to spurious regression a la Granger and Newbold (1974)
- bad identification of random effect vs. serial correlation coefficient is evident (Calzolari and Magazzini 2010)

that can be solved by first-differencing the data. Yet, as Munnell herself (JEL 1992) observes, **first differencing destroys the economic significance of the model.**

So far, the moral of the story is:

- specification analysis of a general model highlighted the limits of our empirical investigation
- panel time series methods might give the correct perspective



Outline of the talk

- 1 A motivating example: the Productivity Puzzle
- 2 Software environment: the **R** project, the **plm** package
- 3 General cross-sectional correlation robustness features
- 4 Diagnostics for *global* or *local* XS dependence
- 5 Digression: testing CSD vs. CWD
- 6 Spatial dependence and spatial models
- 7 ML estimation of spatial panel models
- 8 Software environment: the **splm** package
- 9 Digression: new developments in ML spatial panels
- 10 Answer to the puzzle (?)
- 11 Nonstationary panels

Nonstationarity and XS dependence

The crucial empirical question in a panel time series setting is stationarity

- of variables
- of residuals (cointegration)

Moreover, long panels are prone to common factor heterogeneity, generating

- CSD-type correlation
- unit roots in residuals

Pesaran's augmentation approach accounts for unobserved common factors, allowing

- consistent unit root testing
- consistent and efficient estimation under CSD/CWD

CCE and heterogeneous pts estimators in **plm**

XSD panel time series methods in **plm** (red in R-forge version)

Homogeneous estimators:

- pooled + time FE
- dynamic 2FE (Nickell bias dies out with T)
- **CCEP** (*generalized FE*)

Heterogeneous estimators:

- separate time series regressions (`pvcmm(..., model='within')`)
- Swamy's random coefficients (`pvcmm(..., model='random')`)
- **Mean Groups (MG)**
- **Demeaned MG**
- **CCE MG**



Applied example 5 - CCE pooled and CCEMG

Common Correlated Effects Pooled model

	Estimate	Std. Error	t-value	p-value	
log(pcap)	0.0488771	0.1054583	0.4635	0.6430	
log(pc)	0.0436211	0.0393442	1.1087	0.2676	
log(emp)	0.8376982	0.1415854	5.9166	3.288e-09	***
unemp	-0.0020545	0.0015783	-1.3018	0.1930	

Mean Groups, model=' 'cmg' '

	Estimate	Std. Error	t-value	p-value	
(Intercept)	-0.6741754	1.0445518	-0.6454	0.518655	
log(pcap)	0.0899850	0.1176040	0.7652	0.444180	
log(pc)	0.0335784	0.0423362	0.7931	0.427698	
log(emp)	0.6258659	0.1071719	5.8398	5.225e-09	***
unemp	-0.0031178	0.0014389	-2.1668	0.030249	*



The end

Thanks: in alphabetical order,

- Roger Bivand
- Gianfranco Piras
- Yves, Achim, Christian, Ott, **plm** community
- ...
- ... and you, for your attention