# Spatial Econometrics
Lecture 7: Specifying a spatial model. Case study (1) – wage curve

## Andrzej Torój

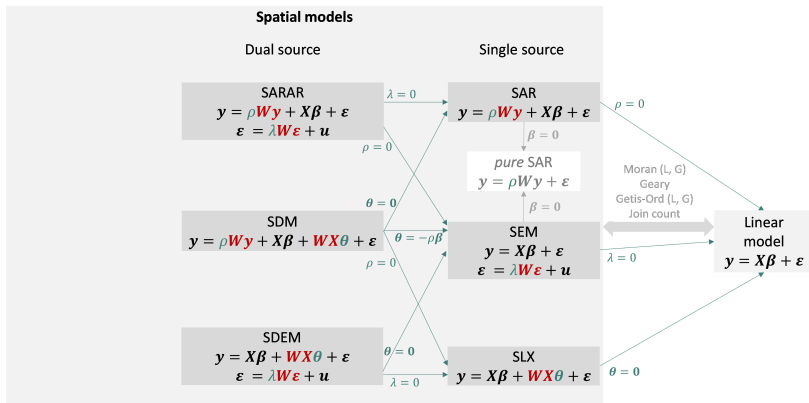Institute of Econometrics – Department of Applied Econometrics

# Outline

1. General Nesting Spatial (GNS) model

2. Strategy of spatial model selection
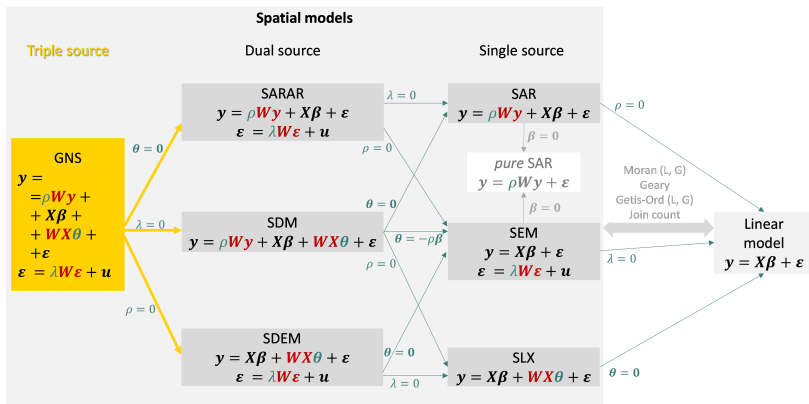
3. Case study: regional wage curve in Poland

# Plan prezentacji

1 General Nesting Spatial (GNS) model

2 Strategy of spatial model selection

3 Case study: regional wage curve in Poland

GNS

# GNS model – specification (1)

# GNS model – specification (2)

# GNS model – specification (3)

- Specification that encompasses all the sources of spatial processes previously under consideration:

$$\mathbf{y} = \rho\mathbf{Wy} + \mathbf{X}\boldsymbol{\beta} + \mathbf{WX}\boldsymbol{\theta} + \boldsymbol{\varepsilon}$$

$$\boldsymbol{\varepsilon} = \lambda\mathbf{W}\boldsymbol{\varepsilon} + \mathbf{u}$$

- Estimation: like SARAR.

# GNS model – specification (3)

- Specification that encompasses all the sources of spatial processes previously under consideration:

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

$$\boldsymbol{\varepsilon} = \lambda\mathbf{W}\boldsymbol{\varepsilon} + \mathbf{u}$$
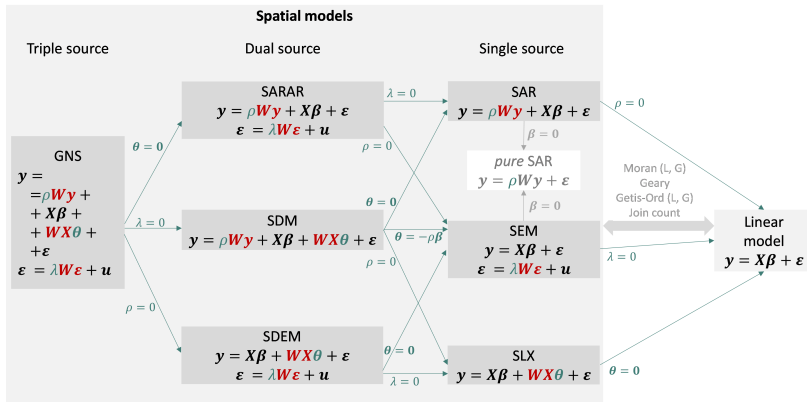
- Estimation: like SARAR.

# Plan prezentacji

1 General Nesting Spatial (GNS) model

2 Strategy of spatial model selection

3 Case study: regional wage curve in Poland

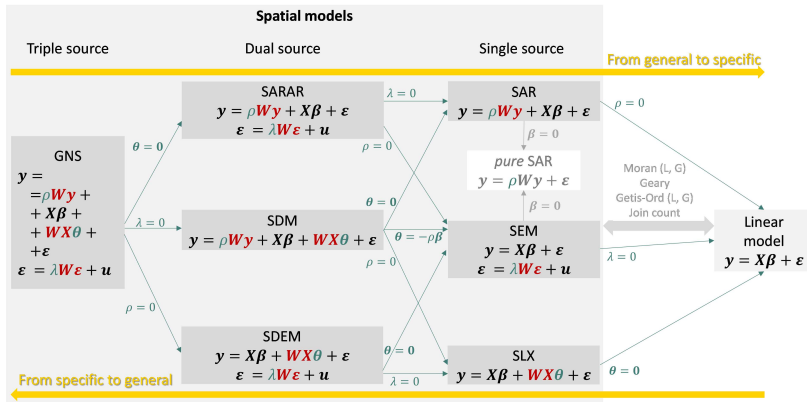GNS model
000

Strategy of spatial model selection
●○○○○○○○○

Case study
○

Selection strategy

# The entire set of models



Scheme inspired by: *Elhorst and Vega (2015)* – with own modifications.

GNS model
○○○

Strategy of spatial model selection
○●○○○○○○○

Case study
○

Selection strategy

# Order of exploration?



Scheme inspired by: *Elhorst and Vega (2015)* – with own modifications.

# From general to specific?

From general to specific:

- Standard approach in time series econometrics (*Hendry, 1995*). We should test from left to right, looking for parsimony in parametrisation and hance higher efficiency.

- Ensures the correct sequential inference as long as no variables are omitted in the general model.

- But GNS is by no means the most general model!

  - One can add further spatial lags ($\mathbf{W}^2$, $\mathbf{W}^3$, ...).
  - Error terms can exhibit local or global spatial autocorrelation.

- Three (or more) sources of spatial process imply **weak statistical identification** of the model (*Cook et al., 2015*):

  - In practice, it is difficult to distinguish between $\rho$, $\boldsymbol{\theta}$ and $\lambda$ (*clustering trilemma*, *label switching problem*.)
  - False decision at the beginning can get the whole inference process side-tracked.

GNS model · · · · · · · · · · · · · · · · · · Strategy of spatial model selection · · · · · · · · · · · · · · · · · · Case study

○○○                                      ○○●○○○○○                                         ○

Selection strategy

# From general to specific?

From general to specific:

- Standard approach in time series econometrics (*Hendry, 1995*). We should test from left to right, looking for parsimony in parametrisation and hance higher efficiency.

- Ensures the correct sequential inference as long as no variables are omitted in the general model.

- But GNS is by no means the most general model!

  - One can add further spatial lags ($\mathbf{W}^2$, $\mathbf{W}^3$, ...).
  - Error terms can exhibit local or global spatial autocorrelation.

- Three (or more) sources of spatial process imply **weak statistical identification** of the model (*Cook et al., 2015*):

  - In practice, it is difficult to distinguish between $\rho$, $\boldsymbol{\theta}$ and $\lambda$ (*clustering trilemma*, *label switching problem*.)
  - False decision at the beginning can get the whole inference process side-tracked.

# From general to specific?

From general to specific:

- Standard approach in time series econometrics (*Hendry, 1995*). We should test from left to right, looking for parsimony in parametrisation and hance higher efficiency.
- Ensures the correct sequential inference as long as no variables are omitted in the general model.
- But GNS is by no means the most general model!
  - One can add further spatial lags ($W^2$, $W^3$, ...).
  - Error terms can exhibit local or global spatial autocorrelation.
- Three (or more) sources of spatial process imply **weak statistical identification** of the model (*Cook et al., 2015*):
  - In practice, it is difficult to distinguish between $\rho$, $\theta$ and $\lambda$ (*clustering trilemma*, *label switching problem*.)
  - False decision at the beginning can get the whole inference process side-tracked.

# From general to specific?

From general to specific:

- Standard approach in time series econometrics (*Hendry, 1995*). We should test from left to right, looking for parsimony in parametrisation and hance higher efficiency.

- Ensures the correct sequential inference as long as no variables are omitted in the general model.

- But GNS is by no means the most general model!

  - One can add further spatial lags ($\mathbf{W}^2$, $\mathbf{W}^3$, ...).
  - Error terms can exhibit local or global spatial autocorrelation.

- Three (or more) sources of spatial process imply **weak statistical identification** of the model (*Cook et al., 2015*):

  - In practice, it is difficult to distinguish between $\rho$, $\boldsymbol{\theta}$ and $\lambda$ (*clustering trilemma*, *label switching problem*.)
  - False decision at the beginning can get the whole inference process side-tracked.

# From specific to general?

From specific to general:

- Less popular a strategy (false partial conclusions due to omitted variable bias).
  - Over-estimating $\beta$ under ommission of the spatial process can make us overlook its presence in the residuals of the linear model.

- Statistical tests frequently suggest many paths "to the left" on the scheme (as well as many paths "to the right" under the previously discussed strategy).

- No testing procedure to directly discriminate between models with identical count of spatial processes (e.g. between single-source SAR vs SEM).

# From specific to general?

From specific to general:

- Less popular a strategy (false partial conclusions due to omitted variable bias).
    - Over-estimating $\beta$ under ommission of the spatial process can make us overlook its presence in the residuals of the linear model.
- Statistical tests frequently suggest many paths "to the left" on the scheme (as well as many paths "to the right" under the previously discussed strategy).
- No testing procedure to directly discriminate between models with identical count of spatial processes (e.g. between single-source SAR vs SEM).
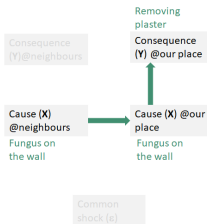
# From specific to general?

From specific to general:

- Less popular a strategy (false partial conclusions due to omitted variable bias).
  - Over-estimating $\beta$ under ommission of the spatial process can make us overlook its presence in the residuals of the linear model.
- Statistical tests frequently suggest many paths "to the left" on the scheme (as well as many paths "to the right" under the previously discussed strategy).
- No testing procedure to directly discriminate between models with identical count of spatial processes (e.g. between single-source SAR vs SEM).
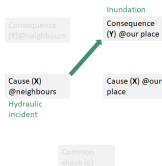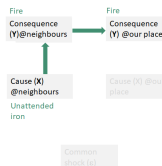
GNS model
○○○

Strategy of spatial model selection
○○○○○●○○○

Case study
○

Selection strategy

# Theory?
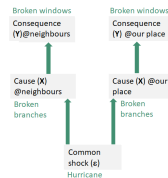
# Get inspired by Box-Jenkins?

In time series econometrics, it holds in stationary and invertible ARMA models that the following conversions are feasible:

- AR(1) into MA($\infty$)
- MA(1) into AR($\infty$)

By the same principle, a limited number of spatial lags of **y** can replace a high number of spatial lags of **X**.
All the tests discussed previously can be helpful.

# Monte Carlo studies?

A very useful guide for empirical researchers written by *Cook et al., (2015)*. Based on Monte Carlo simulations, illustrating the bias in estimation of $\beta, \theta, \rho$ and $\lambda$ under various data generating processes, they conclude that:

- Two recommended solutions are **SARAR** and **SDEM** (unbiased estimation of $\beta$ and maximum chance to avoid spurious conclusions about spatial parameters). When the spatial process is misspecified, the estimates of $\rho$ and $\theta$ from SDM model – heavily biased.

- Researchers that are predominantly interested in unbiased estimates of $\beta$ (but not necessarily in the spatial aspects) should use **SDM**. (*Elhorst, 2011* is much more enthusiastic towards SDM, just like *LeSage* in various publications).

- **GNS** is generally not recommended.

- Single-source models – only when the tests are strongly supportive.

# Monte Carlo studies?

A very useful guide for empirical researchers written by *Cook et al., (2015)*. Based on Monte Carlo simulations, illustrating the bias in estimation of $\beta, \theta, \rho$ and $\lambda$ under various data generating processes, they conclude that:

- Two recommended solutions are **SARAR** and **SDEM** (unbiased estimation of $\beta$ and maximum chance to avoid spurious conclusions about spatial parameters). When the spatial process is misspecified, the estimates of $\rho$ and $\theta$ from SDM model – heavily biased.

- Researchers that are predominantly interested in unbiased estimates of $\beta$ (but not necessarily in the spatial aspects) should use **SDM**. (*Elhorst, 2011* is much more enthusiastic towards SDM, just like *LeSage* in various publications).

- **GNS** is generally not recommended.

- Single-source models – only when the tests are strongly supportive.

# Monte Carlo studies?

A very useful guide for empirical researchers written by *Cook et al., (2015)*. Based on Monte Carlo simulations, illustrating the bias in estimation of $\beta, \theta, \rho$ and $\lambda$ under various data generating processes, they conclude that:

- Two recommended solutions are **SARAR** and **SDEM** (unbiased estimation of $\beta$ and maximum chance to avoid spurious conclusions about spatial parameters). When the spatial process is misspecified, the estimates of $\rho$ and $\theta$ from SDM model – heavily biased.

- Researchers that are predominantly interested in unbiased estimates of $\beta$ (but not necessarily in the spatial aspects) should use **SDM**. (*Elhorst, 2011* is much more enthusiastic towards SDM, just like *LeSage* in various publications).

- **GNS** is generally not recommended.

- Single-source models – only when the tests are strongly supportive.

# Monte Carlo studies?

A very useful guide for empirical researchers written by *Cook et al., (2015)*. Based on Monte Carlo simulations, illustrating the bias in estimation of $\beta, \theta, \rho$ and $\lambda$ under various data generating processes, they conclude that:

- Two recommended solutions are **SARAR** and **SDEM** (unbiased estimation of $\beta$ and maximum chance to avoid spurious conclusions about spatial parameters). When the spatial process is misspecified, the estimates of $\rho$ and $\theta$ from SDM model – heavily biased.

- Researchers that are predominantly interested in unbiased estimates of $\beta$ (but not necessarily in the spatial aspects) should use **SDM**. (*Elhorst, 2011* is much more enthusiastic towards SDM, just like *LeSage* in various publications).

- **GNS** is generally not recommended.

- Single-source models – only when the tests are strongly supportive.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in **LeSage, Pace, 2009**). Nonetheless, they do not give answers to all the doubts listed above.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in **LeSage, Pace, 2009**). Nonetheless, they do not give answers to all the doubts listed above.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in **LeSage, Pace, 2009**). Nonetheless, they do not give answers to all the doubts listed above.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables − intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in **LeSage, Pace, 2009**). Nonetheless, they do not give answers to all the doubts listed above.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in *LeSage, Pace, 2009*). Nonetheless, they do not give answers to all the doubts listed above.

GNS model
000

Strategy of spatial model selection
○○○○○○○●

Case study
○

Selection strategy

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in **LeSage, Pace, 2009**). Nonetheless, they do not give answers to all the doubts listed above.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in **LeSage, Pace, 2009**). Nonetheless, they do not give answers to all the doubts listed above.

# Choosing the estimation method

When the specification is determined, the decision is much easier here. For almost all types of models, we can choose between ML-type estimation and LS/GMM-type estimation.

- In standard cases: no recommendations.
- Large $N$ (e.g. a few thousand): ML estimators evaluate the determinant of $N \times N$ matrix at each iteration. A lot of time and computational capacity is needed.
- LS / GMM methods are much more convenient when spatial autocorrelation is accompanied by **endogeneity of other regressors** (adding further instrumental variables – intuitive).
- ML framework implies very intuitive tools for testing nested models (e.g. test LR).
- ML sometimes ensures higher efficiency, as it consumes the knowledge of entire error distribution, instead of only error moments. We cannot be sure, however, that the usual assumption of error normality is fulfilled.
- LS / GMM methods more robust to specification errors.
- Many authors use Bayesian methods with spatial models (cf. thorough discussion in *LeSage, Pace, 2009*). Nonetheless, they do not give answers to all the doubts listed above.

## Plan prezentacji

1 General Nesting Spatial (GNS) model

2 Strategy of spatial model selection

3 Case study: regional wage curve in Poland

# Exercise

Using all the previously studied techniques and the dataset used in lecture 1 (for Polish poviats), investigate the impact of unemployment rate on the wages in Polish poviats.

1. Determine wheter spatial modelling is relevant here. Construct the weight matrix (look at 2-3 alternative options).

2. Find the best model.

3. Add further variables to the equation if necessary.

4. Select one poviat (but not a city-poviat) and, based on the obtained model, illustrate how a location of a foreign greenfield investment will affect the wage levels in various poviats. Assume that the investment itself will directly decrease the unemployment rate in this poviat by 0.5 p.p.