# Parallel PySAL

## Autoregression and Complex Systems Framework Integration

Jason Laura, Robert Pahle, Sergio Rey, Luc Anselin

GeoDa Center for Geospatial Analysis and Computation
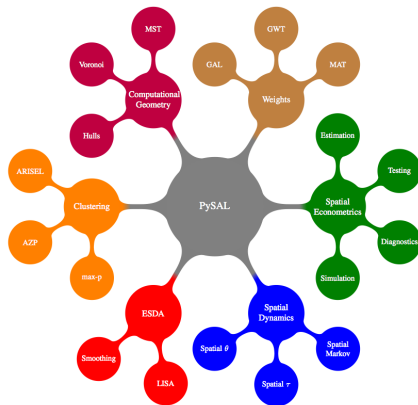Arizona State University

August 18, 2014

# Outline

PySAL

Substantive Application: Spatial Econometrics
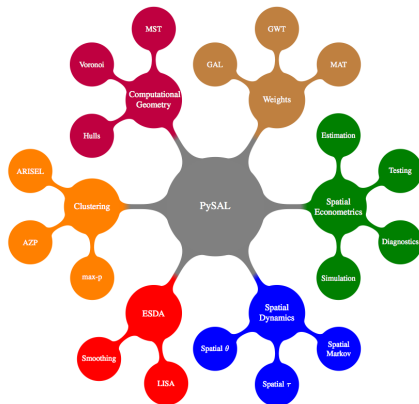
Implementation

# PySAL

- Spatial analysis library
- Big data world
- v 1.8 July 2014

# pPySAL

- contiguity builder
- max-p region
- p-lisa
- fisher jenks
- spatial regimes

## Lessons Learned

- ▶ Hardware dependence
- ▶ No holy grail of automatic parallelization
- ▶ Need a roadmap = Taxonomy
  - ▶ Guidance on "best practice"
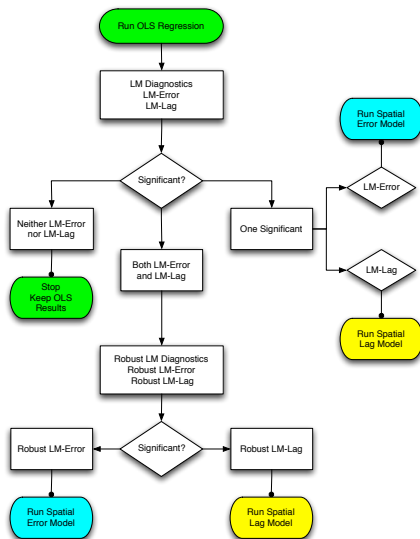  - ▶ Identify dead ends

# GeoDaSpace: Spatial Econometrics

- ▶ GUI ontop of spreg
- ▶ Subset of spreg functionality
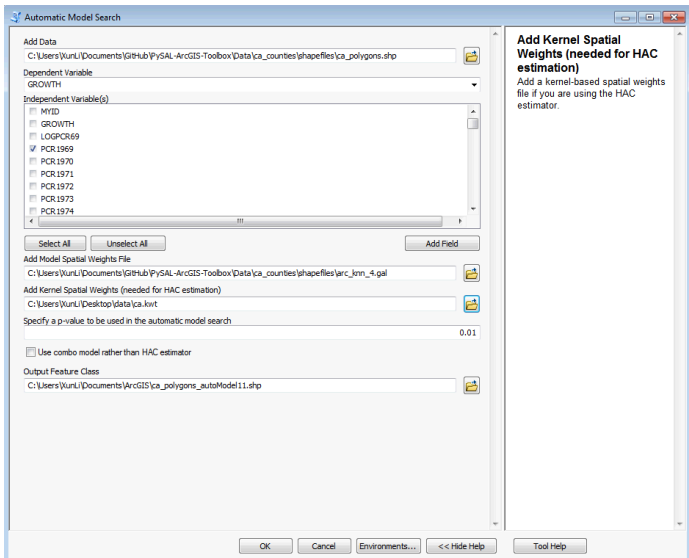- ▶ Cross-platform

## Specification Searches

- ▶ Specific to General
    - ▶ $y = X\beta + \epsilon$
    - ▶ OLS + Lagrange Multiplier Tests
- ▶ General to Specific
    - ▶ $y = \rho Wy + X\beta + (I - \lambda W)^{-1}\nu$
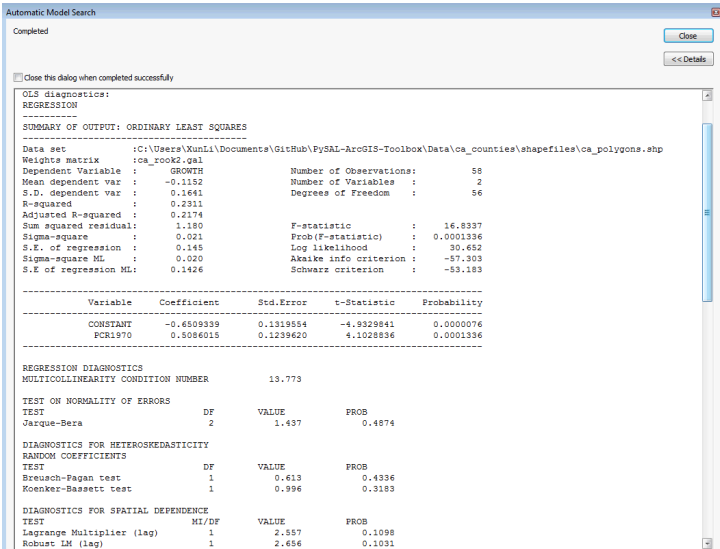    - ▶ ML + Restrictions

# LM Based Specification

# ArcGIS Toolbox

# ArcGIS Toolbox

Root Node: Ordinary Least Squares Regression

Then

A. If Lagrange Multiplier Test for Spatial Error Model < p-value AND Lagrange Multiplier Test for Spatial Lag Model < p-value

    1. If Robust Lagrange Multiplier Test for Spatial Error p-value < p-value and Robrust Lagrange Multiplier Test for Spatial Lag Model p-value < p-vlaue:
        a. If NOT combo
          i. twosls_sp.GM_Lag
          ii. "Spatial Lag with Spatial Error - HAC"
        b. Elif Koenker Basset Statistic p-value < p-value
          i. error_sp_het.GM_Combo_Het
          ii. "Spatial Lag with Spatial Error - Heteroskedastic"
        c. Else
          i. error_sp_hom.GM_Combo_Hom
          ii. "Spatial Lag with Spatial Error - Homoskedastic"

    2. Else If Robust Lagrange Multiplier Test for Spatial Error p-value < p-value and RLM for Spatial Lag p-value > p-value:
        a. If OLS Koenker Basset Statistic p-value < p-value
          i. error_sp_het.GM_Error_Het
          ii. "Spatial Error - Heteroskedastic"
        b. Else If OLS Koenker Basset Statistic p-value > p-value
          i. error_sp_hom.GM_Error_Hom
          ii. "Spatial Error - Homoskedastic"

    3. Else If RLM for Spatial Error > p-value and RLM for Spatial Lag < p-value
        a. If OLS Koenker Basset Statistic p-value < p-value
          i. twosls_sp.GM_Lag (robut:white)
          ii. "Spatial Lag - Heteroskedastic"
        b. Else If OLS Koenker Basset Statistic p-value > p-value
          i. twosls_sp.GM_Lag
          ii. "Spatial Lag - Homoskedastic"
    4. Else If RLM for Spatial Error > p-value and RLM for Spatial Lag > p-value
        a. No PySAL Call
        b. No Model - Robust Test not Significant - Check Model.

B. Else If Lagrange Mutiplier Test for Spatial Error Model < p-value AND Lagrange Multiplier Test for Spatial Lag > p-value
    1. If OLS Koenker Basset Statistic p-value < p-value
        i. error_sp_het.GM_Error_Het
        ii. "Spatial Error - Heteroskedastic"
    2. Else If OLS Koenker Basset Statistic p-value > p-value
        i. error_sp_hom.GM_Error_Hom
        ii. "Spatial Error - Homoskedastic"

C. Else If Lagrange Multiplier Test for Spatial Error Model > p-value AND Lagrange Multiplier Test for Spatial Lag < p-value
    1. If OLS Koenker Basset Statistic p-value < p-value
        i. twosls_sp.GM_Lag (robust-white)
        ii. "Spatial Lag - Heteroskedastic"
    2. Else If OLS Koenker Basset Statistic p-value > p-value
        i. twosls_sp.GM_Lag
        ii. "Spatial Lag - Homoskedastic"

D. Else Lagrange Multiplier Test for Spatial Error Model > p-value AND Lagrange Multiplier Test for Spatial Lag > p-value

    1. If OLS Koenker Basset Statistic p-value < p-value
        i. ols.OLS (robust-white)
        ii. "No Space - Heteroskedastic"
    2. Else If OLS Koenker Basset Statistic p-value > p-value
        i. ols.OLS
        ii. "No Space - Homoskedastic"

# Parallel Strategy

- ▶ Speculative Parallelism
    - ▶ Solve' all branches of a search tree
    - ▶ Leverage an excess computation model
    - ▶ No dependency in execution order
    - ▶ Synchronization at the completion of all computation
- ▶ Implementation
    - ▶ Utilize a processing queue
    - ▶ One manager, and n workers
    - ▶ Workers draw a regression model from the queue, process, and return the result
    - ▶ Scales to where n = number of models to compute
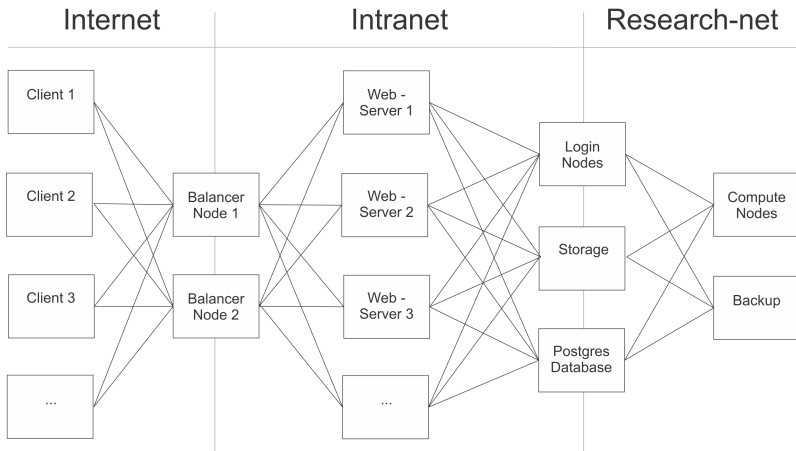    - ▶ Potential to extend to variable parameter specification (larger tree)

# Tensions

## Trade-off

- ▶ Trading elegant econometric theory for data mining
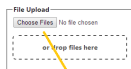- ▶ Gain speed and coverage of model space over the sequential approach

## Issues

- ▶ Distributional properties of big data approach unknown
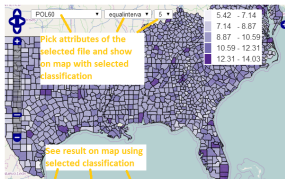- ▶ Purists take a dim view of "data mining"

# Complex Systems Framework

# Autoreg in CSF

# Model Path



**Autoreg decision tree**

1. Robust Lagrange Multiplier Test for:
1) Spatial Error p-value < p-value and
2) Spatial Lag Model p-value < p-value

NOT combo:
->twosls_sp GM_Lag
->Spatial Lag with Spatial Error - HAC

Koenker Basset Statistic p-value < p-value
->error_sp_hat GM_Combo_Hat
->Spatial Lag with Spatial Error - Heteroskedastic

Else
->error_sp_hom GM_Combo_Hom
->Spatial Lag with Spatial Error - Homoskedastic

2. Robust Lagrange Multiplier Test for:
1) Spatial Error p-value < p-value and
2) Spatial Lag p-value > p-value

OLS Koenker Basset Statistic p-value < p-value
->error_sp_hat GM_Error_Hat
->Spatial Error - Heteroskedastic

OLS Koenker Basset Statistic p-value > p-value
->error_sp_hom GM_Error_Hom
->Spatial Error - Homoskedastic

A. Lagrange Multiplier Test for:
1) Spatial Error Model < p-value and
2) Spatial Lag Model < p-value

3. Robust Lagrange Multiplier Test for:
1) Spatial Error p-value > p-value and
2) Spatial Lag p-value < p-value

OLS Koenker Basset Statistic p-value < p-value
->twosls_sp GM_Lag (robust-white)
->Spatial Lag - Heteroskedastic

OLS Koenker Basset Statistic p-value > p-value
->twosls_sp GM_Lag
->Spatial Lag - Homoskedastic

4. Robust Lagrange Multiplier Test for:
1) Spatial Error p-value > p-value and
2) Spatial Lag p-value > p-value

No PySAL Call

No Model
Robust Test not Significant
Check Model

Ordinary Least
Squares Regression

B. Lagrange Multiplier Test for:
1) Spatial Error Model < p-value and
2) Spatial Lag Model > p-value

OLS Koenker Basset Statistic p-value < p-value
->error_sp_hat GM_Error_Hat
->Spatial Error - Heteroskedastic

OLS Koenker Basset Statistic p-value > p-value
->error_sp_hom GM_Error_Hom
->Spatial Error - Homoskedastic

**Selected path for the current result**

C. Lagrange Multiplier Test for:
1) Spatial Error Model > p-value and
2) Spatial Lag Model < p-value

OLS Koenker Basset Statistic p-value < p-value
->twosls_sp GM_Lag (robust-white)
->Spatial Lag - Heteroskedastic

OLS Koenker Basset Statistic p-value > p-value
->twosls_sp GM_Lag
->Spatial Lag - Homoskedastic

OLS Koenker Basset Statistic p-value < p-value

# Next Steps

## Parallel Autoreg

- ► Ensemble of search strategies
  - ► short
  - ► full
  - ► hybrid
- ► Candidate Variables
- ► Candidate $W$s

## Integration

- ► CyberGIS Gateway
- ► Strategies

Come see the demo!