

REGRESSION MODELS FOR SPATIALLY AUTOCORRELATED DATA



Meïli Baragatti

OpenSpat

Table of Contents

1 Origins and Consequences of Spatial Autocorrelation

- Origin : interaction
- Origin : reaction
- Origin : misspecification
- Consequences of the spatial autocorrelation on classical linear models

2 Working example : Las Rosas

3 Spatial Lag Model

4 Spatial Error Model

5 Choosing Between Spatial Lag, Error and SAC models

6 Extended Linear Models

7 Bibliography

Origins and Consequences of Spatial Autocorrelation

When we detect an apparent spatial autocorrelation (on residuals for instance), this spatial autocorrelation may or may not be the result of a spatial autocorrelation.

In 1984, Miron identified three Origins of apparent or real spatial autocorrelation :

- interaction
- reaction
- misspecification

Origins of spatial autocorrelation : example

Imagine a population of plants growing in a particular region :

- Y_i measurement of plant productivity (tree height or population density).
- Population is sufficiently dense relative to the spatial scale \Rightarrow productivity measurement may be modeled as varying continuously with the location.
- X_{i1} the amount of light available at location i .
- X_{i2} the amount of available nutrients at location i .

Using these two explanatory variables, the simplest model is :

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i, \quad \text{with } \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2). \quad (1)$$

In matrix notation :

$$\begin{aligned} Y &= X\beta + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I). \end{aligned} \quad (2)$$

The following three notions can be combined in a same model.

Origin : interaction

Spatial autocorrelation induced by interaction occurs when the response variables at different sites interact with each other.

- **Negative autocorrelation** may occur if trees in close proximity compete with each other for light and nutrients, so that relatively productive tree populations tend to inhibit the growth of other trees.
- **Positive autocorrelation** would occur if existing trees produced acorns that do not disperse very far, which in turn results in more trees in the vicinity.

If Y is positively autocorrelated, the true underlying model is :

$$\begin{aligned} Y &= X\beta + \rho WY + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I), \end{aligned} \tag{3}$$

with WY the spatial lag.

Interaction : illustration using simulations (1/2)

We generate a dataset `simu_modlin` satisfying model (2) with $\beta = (0, 0.5, 0.3)$ and a dataset `simu_interaction` satisfying model (3) with $\beta = (0, 0.5, 0.3)$ and $\rho = 0.6$. Each dataset contains 1000 observations and X_1 and X_2 are simulated independently using gaussian distributions.

```
mod <- lm(Ylin ~ X1 + X2)
print(coef(mod), digits = 2)

## (Intercept)           X1           X2
## -0.00021      0.49979      0.30028

var(mod$res)

## [1] 9.560601e-05
```

Interaction : illustration using simulations (2/2)

```
mod <- lm(Yinter ~ X1 + X2)
print(coef(mod), digits = 2)

## (Intercept)          X1          X2
##      0.027       0.550       0.327

var(mod$res)

## [1] 0.06299029
```

Origin : reaction

Spatial autocorrelation induced by reaction occurs when the response variables are reacting to an external factor that varies in space, and when this factor is not taken into account by the model.

For instance if nearby plants are reacting to availability of water (which varies in the 'space').

The inclusion of this external factor in the linear model may be appropriate. It may be sufficient to explain the spatial autocorrelation, and to obtain non-autocorrelated residuals.

For instance, the true model should be :

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \epsilon_i, \quad \text{with } \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2), \quad (4)$$

with X_{i3} the distance from the river at location i .

Reaction : illustration using simulations (1/4)

We generate a dataset `simu_reaction1` satisfying model (4) with $\beta = (0, 0.5, 0.3, 0.8)$ and X_3 correlated with X_2 .

We fit model (4) :

```
print(coef(lm(Yreact1 ~ X1 + X2 + X3)), digits = 2)
```

```
## (Intercept)          X1          X2          X3  
##     0.0088     0.4837     0.3315     0.7716
```

```
mod <- lm(Yreact1 ~ X1 + X2)  
print(coef(mod), digits = 2)
```

```
## (Intercept)          X1          X2  
##     0.51        0.50        1.01
```

X_3 maybe interpreted as a ‘spatial’ variable, but its role in the model is identical to that of another explanatory variable without any spatial connotation.

Reaction : illustration using simulations (2/4)

We generate a dataset `simu_reaction1` satisfying model (4) with $\beta = (0, 0.5, 0.3, 0.8)$, and X_3 non correlated with X_1 or X_2 but spatially autocorrelated.

We fit model (4) :

```
mod <- lm(Yreact2 ~ X1 + X2 + X3)
print(coef(mod), digits = 2)

## (Intercept)          X1          X2          X3
##      0.027       0.483       0.327       0.777

var(mod$res)

## [1] 1.004552
```

Reaction : illustration using simulations (3/4)

We fit model (1) :

```
mod <- lm(Yreact2 ~ X1 + X2)
print(coef(mod), digits = 2)

## (Intercept)          X1          X2
##      0.051       0.469       0.312

var(mod$res)

## [1] 1.90547
```

The effect of X_3 which is not taken into account in this model is entirely loaded in the error term.

Reaction : illustration using simulations (4/4)

As X_3 was spatially autocorrelated, the result is that the residuals are spatially autocorrelated :

```
lm.morantest(mod,W)

##
## Global Moran I for regression residuals
##
## data:
## model: lm(formula = Yreact2 ~ X1 + X2)
## weights: W
##
## Moran I statistic standard deviate = 6.7573, p-value = 7.027e-12
## alternative hypothesis: greater
## sample estimates:
## Observed Moran I      Expectation      Variance
##          0.1550852149 -0.0010125073  0.0005336302
```

Origin : misspecification

The measured autocorrelation is not due to interaction or reaction but to the incorrect form of the model.

For instance if we assume homoscedastic errors when in fact they are heteroscedastic.

The true model should be (the variance of the errors increases with the amount of available nutrients X_{i2}) :

$$\begin{aligned} Y &= X\beta + \epsilon \\ \epsilon_i &\stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2 \times \exp(1 + 2X_{i2})). \end{aligned} \tag{5}$$

In this case, the measured autocorrelation can be induced by the wrong modelisation, it is then an apparent autocorrelation and not a real autocorrelation (this autocorrelation cannot be explained by spatial considerations).

Misspecification : illustration using simulations (1/2)

We generate a dataset `simu_modmiss` satisfying model (5) with $\beta = (0, 0.5, 0.3)$. X_2 spatially autocorrelated and the error variance is an increasing function of X_2 .

We fit model (2) :

```
mod <- lm(Ymiss ~ X1 + X2)
print(coef(mod, digits = 2))

## (Intercept)          X1          X2
##   30.95456  -68.36515   84.34902
```

Misspecification : illustration using simulations (2/2)

```
lm.morantest(mod, W)

##
## Global Moran I for regression residuals
##
## data:
## model: lm(formula = Ymiss ~ X1 + X2)
## weights: W
##
## Moran I statistic standard deviate = 2.3661, p-value = 0.008989
## alternative hypothesis: greater
## sample estimates:
## Observed Moran I      Expectation      Variance
##      0.159988117     -0.014521469     0.005439884
```

The error terms are uncorrelated, but because the error variance is a function of X_2 and high values of X_2 tend to be near other high values of X_2 , a test for spatial autocorrelation of the residuals has a high type I error rate.

Consequences of the spatial autocorrelation (1/2)

Whatever the origin of apparent two-dimensional spatial autocorrelation, effects of this autocorrelation similar to those for one-dimensional autocorrelation : autocorrelation decreases the effective sample size, as there are no longer n independent sources of information.

⇒ the standard statistical techniques which are derived under the assumption of independence will provide mistaken significance levels and p -values, as well as mistaken confidence levels for confidence intervals.

Consequences of the spatial autocorrelation (2/2)

- Interaction** biased estimates of the regression coefficients, the variance of the residuals is inflated \Rightarrow inflated type I or II error rates of certain tests.
- reaction** If the reaction variable (not included in the model) is correlated to a variable present in the model, the estimate of the coefficient associated with the variable present in the model will be biased.
If the reaction variable (not included in the model) is not correlated to a variable present in the model, but is spatially autocorrelated, the variance of the residuals will be inflated,
 \Rightarrow inflated type I or II error rates and indication of spatial autocorrelation when none really exists.
- Misspecification** If the model is misspecified, that can lead to both biased estimates of the regression coefficient and indication of spatial autocorrelation when none really exists.

Table of Contents

- 1 Origins and Consequences of Spatial Autocorrelation
- 2 Working example : Las Rosas
- 3 Spatial Lag Model
- 4 Spatial Error Model
- 5 Choosing Between Spatial Lag, Error and SAC models
- 6 Extended Linear Models
- 7 Bibliography

Spatial regression models in practice (1/2)

1 Fit the data with a classical linear model like (2).

2 Check the model assumptions on the residuals : normality, homoscedasticity and independence.

Non-normality histogram, Q-Q plot.

Heteroscedasticity or the exclusion of a reaction variable plot the residuals against the fitted values, and against the different variables included or not in the model.

Dependence try to detect a spatial autocorrelation of the residuals : bubble plots, semi-variograms, Moran correlogram, test for spatial autocorrelation of the residuals using the Moran's I .

Spatial regression models in practice (2/2)

3 If we detect some problems on the residuals :

Non-normality the model can be misspecified. Try a transformation of your variable to be explained and/or of your explanatory variables. It can also be the consequence of a relevant explanatory variable forgotten in the model.

Homoscedasticity or the exclusion of a reaction variable take into account this heteroscedasticity in your model.

Dependence check that you have not forgotten a reaction variable, and that you are not in presence of heteroscedasticity. If not, fit a more complicated model with an autocorrelation structure : spatial lag model, spatial error model or an extended linear model with a spatial autocorrelation structure.

Example Las Rosas (1/23)

Data set from Anselin et al. 2001.

- Measurements of **corn yield** over a controlled plot in Argentina. Regular grid approximately 71 cm apart.
- **Amount of nitrogen fertilizer** that is applied on each location : 6 levels applied along the rows of the field.
- The basic set of information consists of four variables measured at 1704 locations : YIELD, N, LATITUDE, LONGITUDE.
- Xutm a SpatialPointsDataFrame object containing the yield and relevant geographical variables to explain it (N, elev, slope, slopeX, accu, aspect and hshade).

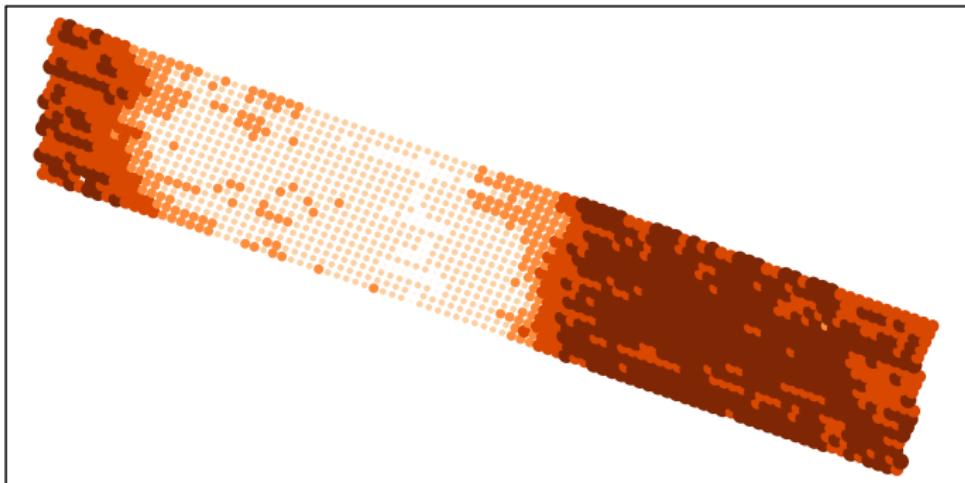
Objective

Do some of the explanatory variables influenced the observed yield variability in the field ?

Example Las Rosas (2/23)

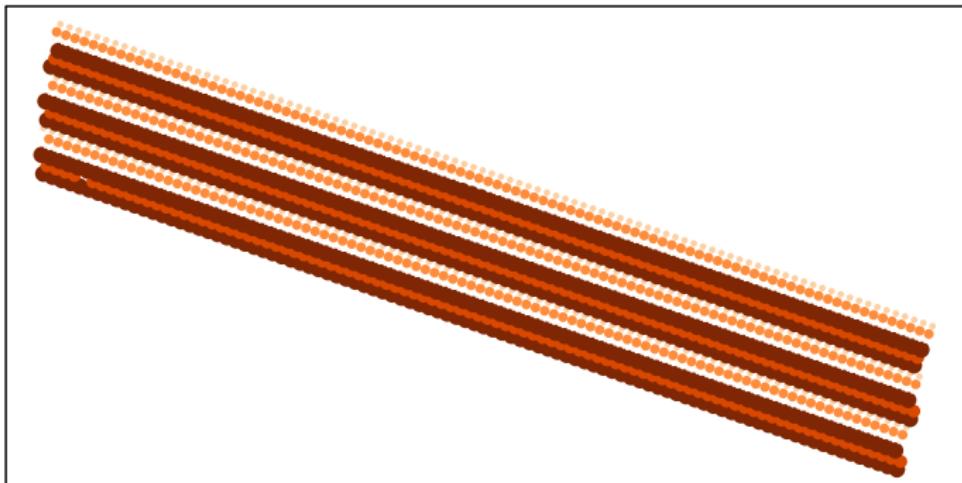
```
## Formal class 'SpatialPointsDataFrame' [package "sp"] with 5 slots
##   ..@ data      : 'data.frame': 1704 obs. of 10 variables:
##     ...$ YIELD : num [1:1704] 4225 4308 4301 4443 4343 ...
##     ...$ N      : num [1:1704] 125 125 125 125 125 ...
##     ...$ elev   : num [1:1704] 272 272 272 272 272 ...
##     ...$ slope  : num [1:1704] 0.022 0.0238 0.0256 0.027 0.0282 ...
##     ...$ slopeX: num [1:1704] 13.4 14.5 15.7 16.6 17.4 ...
##     ...$ accu   : num [1:1704] 72.5 70.4 68.3 65.9 61.2 ...
##     ...$ aspect  : num [1:1704] 4.46 4.5 4.53 4.55 4.58 ...
##     ...$ hshade : num [1:1704] 0.864 0.864 0.864 0.864 0.864 ...
##     ...$ x      : num [1:1704] 420774 420781 420787 420794 420800 ...
##     ...$ y      : num [1:1704] 6342855 6342853 6342850 6342847 6342845 ...
##   ..@ coords.nrs : num(0)
##   ..@ coords    : num [1:1704, 1:2] 420774 420781 420787 420794 420800 ...
##   ... - attr(*, "dimnames")=List of 2
##     ... .$. : chr [1:1704] "1" "2" "3" "4" ...
##     ... .$. : chr [1:2] "LONGITUDE" "LATITUDE"
##   ..@ proj4string:Formal class 'CRS' [package "sp"] with 1 slot
##   ... .@ projargs: chr "+proj=utm +zone=20 +south +ellps=WGS84 +datum=SP
```

Example Las Rosas (3/23)

Yield

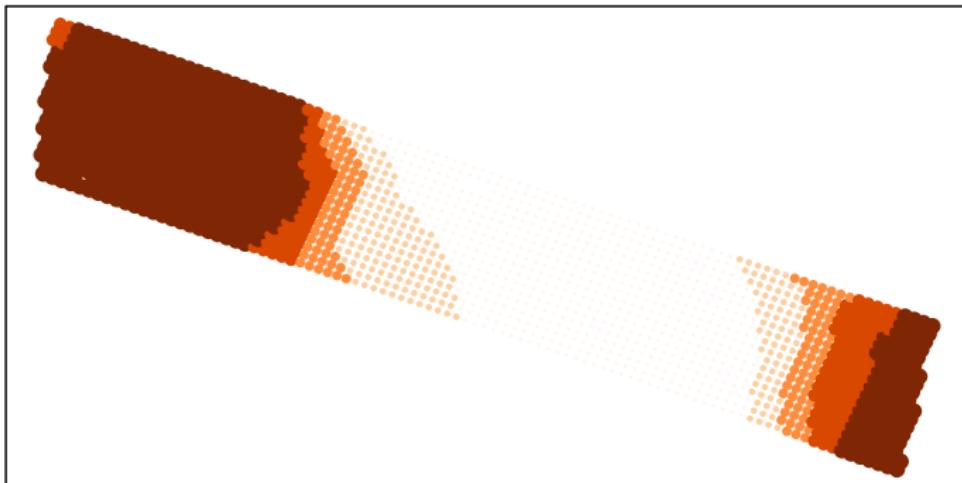
- [574.2,1529]
• (1529,2484)
• (2484,3438)
• (3438,4393)
• (4393,5348)

Example Las Rosas (4/23)

Nitrogen

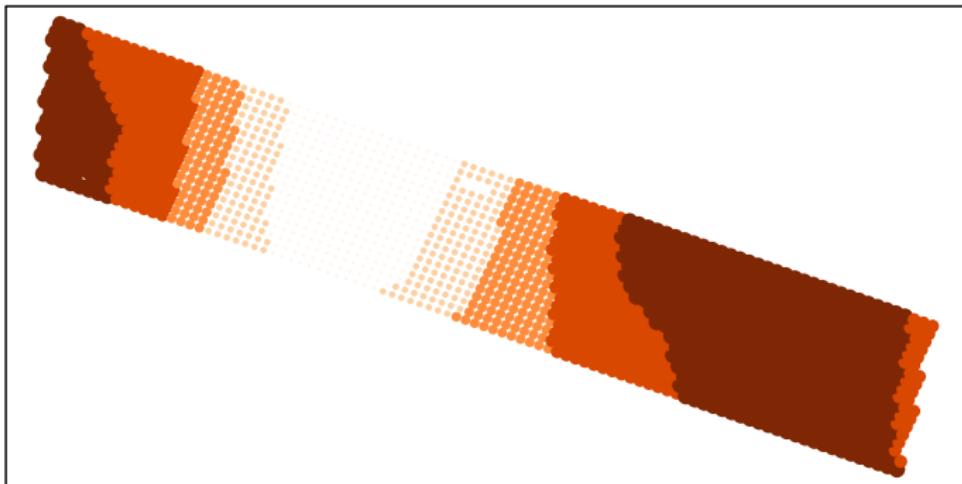
- [0,24.92]
- (24.92,49.84]
- (49.84,74.76]
- (74.76,99.68]
- (99.68,124.6]

Example Las Rosas (5/23)

Soil aspect

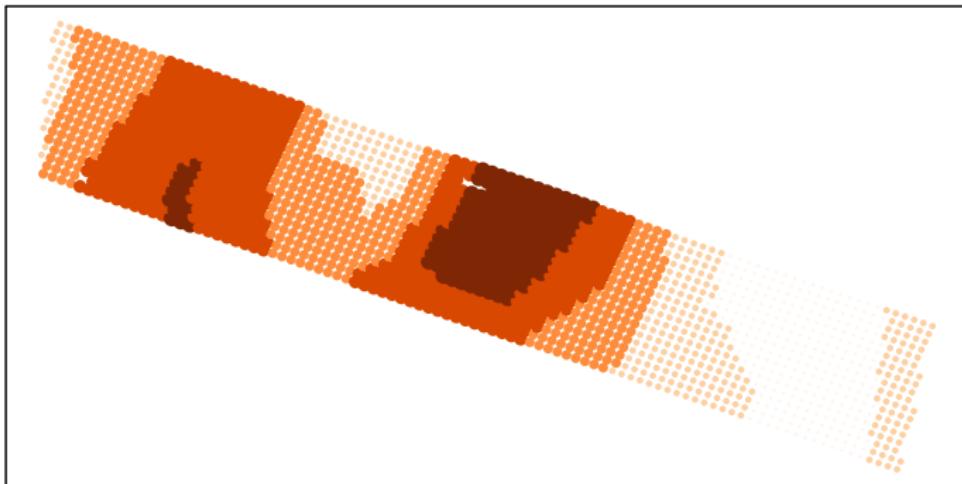
- [1.85,2.439]
- (2.439,3.027]
- (3.027,3.615]
- (3.615,4.203]
- (4.203,4.792]

Example Las Rosas (6/23)

Water accumulation

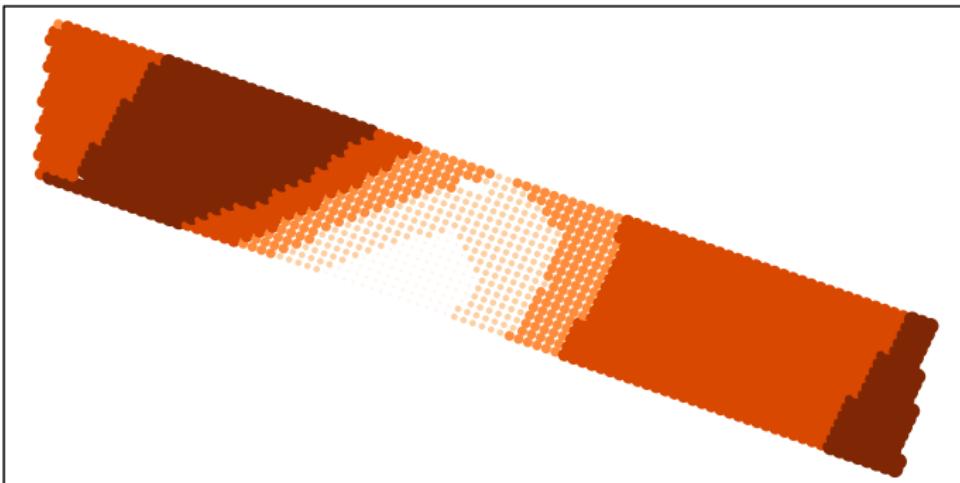
- [-143.1,-96.59]
- (-96.59,-50.04]
- (-50.04,-3.495]
- (-3.495,43.05]
- (43.05,89.6]

Example Las Rosas (7/23)

Slope

- [0.007235,0.0141]
- (0.0141,0.02097]
- (0.02097,0.02783]
- (0.02783,0.0347]
- (0.0347,0.04156]

Example Las Rosas (8/23)

Amount of radiation

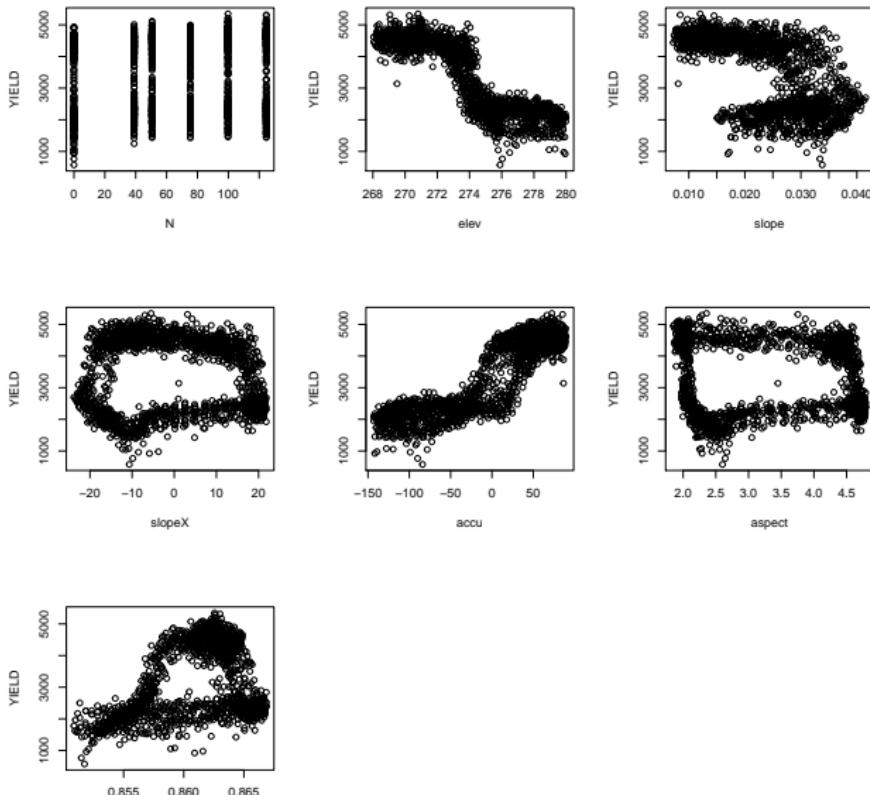
- [0.8509,0.8541]
- (0.8541,0.8573]
- (0.8573,0.8605]
- (0.8605,0.8637]
- (0.8637,0.8669]

Example Las Rosas (9/23)

```
round(cor(Xutm@data[,1:8]),3)
```

```
##          YIELD         N      elev    slope slopeX     accu aspect hshade
## YIELD  1.000  0.079 -0.881 -0.627 -0.107  0.889 -0.144  0.378
## N      0.079  1.000 -0.022  0.008  0.003 -0.001  0.002 -0.043
## elev   -0.881 -0.022  1.000  0.584  0.123 -0.954  0.108 -0.306
## slope  -0.627  0.008  0.584  1.000 -0.051 -0.525  0.033 -0.368
## slopeX -0.107  0.003  0.123 -0.051  1.000  0.016  0.965  0.708
## accu   0.889 -0.001 -0.954 -0.525  0.016  1.000  0.015  0.424
## aspect  -0.144  0.002  0.108  0.033  0.965  0.015  1.000  0.613
## hshade  0.378 -0.043 -0.306 -0.368  0.708  0.424  0.613  1.000
```

Example Las Rosas (10/23)



Example Las Rosas (11/23)

- Relationships between the yield and elev and between the yield and accu are similar.
- Relationships between the yield and aspect and between the yield and slopeX are similar.
- Some scatterplots can be separated in two scatterplots having different relationships with the yield.
- We suspect that the separation can be made regarding on the elevation, or the water accumulation.

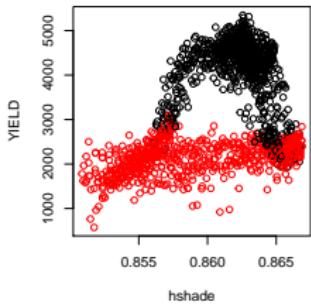
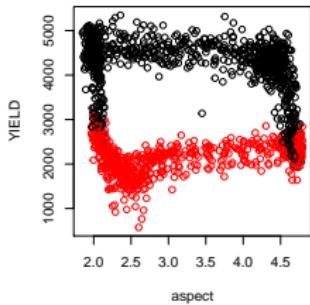
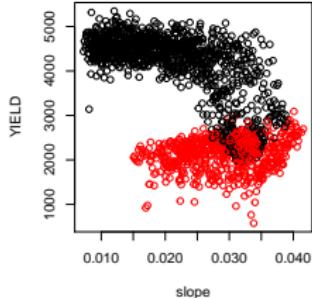
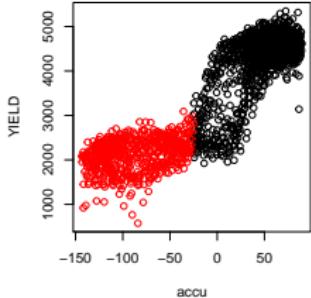
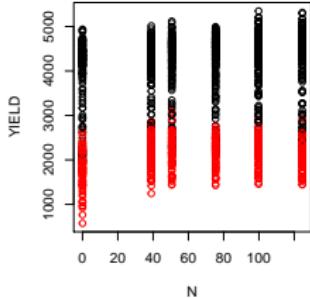
We decide to keep only the water accumulation in our model, and to transform it into a factor with two levels : low and high.

Example Las Rosas (12/23)

We transform the variable accu into a factor with two levels : low and high :

```
Xutm@data$accuf <- rep('low',dim(Xutm@data)[1])
for (i in 1:dim(Xutm@data)[1]){
  if(Xutm@data$accu[i] > -25){Xutm@data$accuf[i] <- 'high'}
}
Xutm@data$accuf <- as.factor(Xutm@data$accuf)
```

Example Las Rosas (13/23)



Example Las Rosas (14/23)

A linear regression model `model.lm` is proposed to explained the yield using all these explanatory variables and their interactions

```
model.lm <- lm(YIELD ~ accuf + N + slope + aspect + hshade + accuf*slope
                 + accuf*aspect + accuf*hshade, data=Xutm@data)
drop1(model.lm,test="F")
```

```
## Single term deletions
##                   Df Sum of Sq      RSS      AIC F value    Pr(>F)
## <none>              276691762 20462
## N                  1   18869062 295560825 20572 115.591 < 2.2e-16 ***
## accuf:slope        1   96569568 373261331 20970 591.580 < 2.2e-16 ***
## accuf:aspect        1   3455382 280147145 20481 21.168 4.521e-06 ***
## accuf:hshade        1   16482486 293174248 20559 100.971 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Example Las Rosas (15/23)

```
summary(model.lm)
```

```
## Coefficients:  
##                                     Estimate Std. Error t value Pr(>|t|)  
## (Intercept)      6.570e+04  8.007e+03   8.204 4.53e-16 ***  
## accuflow       -1.004e+05  9.520e+03  -10.549 < 2e-16 ***  
## N                2.601e+00  2.419e-01   10.751 < 2e-16 ***  
## slope          -5.925e+04  1.626e+03  -36.436 < 2e-16 ***  
## aspect          -1.600e+02  1.643e+01   -9.739 < 2e-16 ***  
## hshade          -6.949e+04  9.314e+03  -7.461 1.37e-13 ***  
## accuflow:slope  7.863e+04  3.233e+03   24.322 < 2e-16 ***  
## accuflow:aspect 1.471e+02  3.197e+01    4.601 4.52e-06 ***  
## accuflow:hshade 1.116e+05  1.110e+04   10.048 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 404 on 1695 degrees of freedom  
## Multiple R-squared:  0.8796, Adjusted R-squared:  0.879  
## F-statistic: 1548 on 8 and 1695 DF,  p-value: < 2.2e-16
```

Example Las Rosas (16/23)

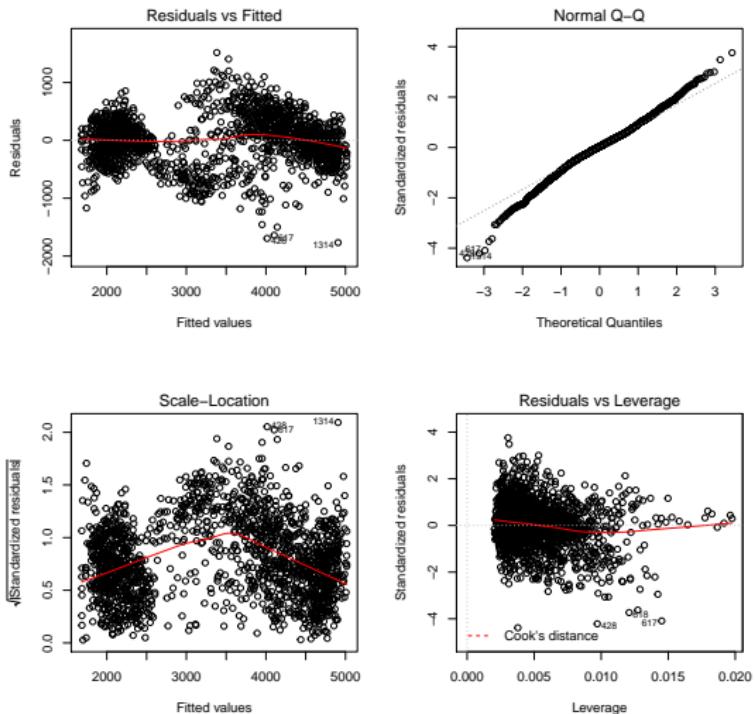
The equation of the model is the following, with $Yield_{ij}$ the value of yield at the j^{th} location having level i of accu.

$$Yield_{ij} = \beta_0 + \alpha_i + \beta_1 N_{ij} + \beta_2 slope_{ij} + \beta_3 aspect_{ij} + \beta_4 hshade_{ij} + \gamma_2 i slope_{ij} + \gamma_3 i aspect_{ij} + \gamma_4 i hshade_{ij} + \epsilon_{ij}, \quad (6)$$

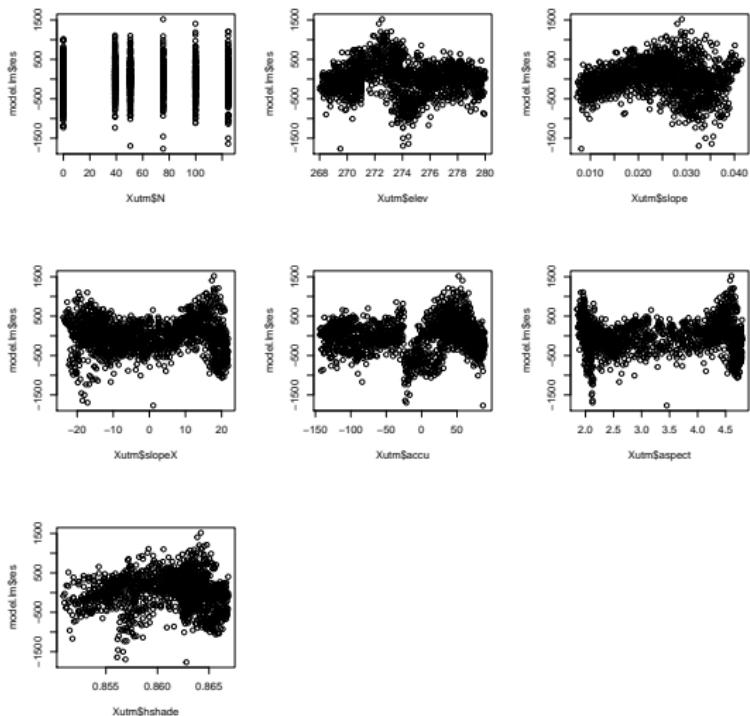
$$\epsilon_{ij} \underset{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2). \quad (7)$$

We then need to validate the assumptions.

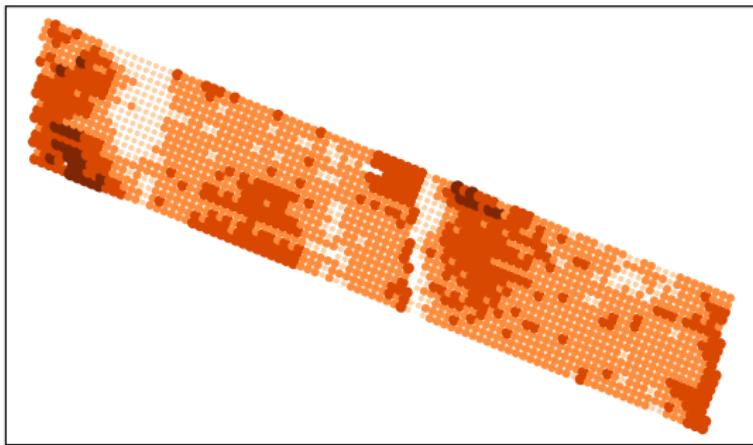
Example Las Rosas (17/23)

Figure 1 – Diagnostic plots for `model.lm`.

Example Las Rosas (18/23)

Figure 2 – Residuals of `model.lm` against every possible explanatory variable.

Example Las Rosas (19/23)

Residuals of model.lm

- $[-1767, -1110]$
- $\{-1110, -453.4\}$
- $\{-453.4, 203.5\}$
- $(203.5, 860.4)$
- $(860.4, 1517]$

Figure 3 – Bubble map for residuals of model.lm.

Example Las Rosas (20/23)

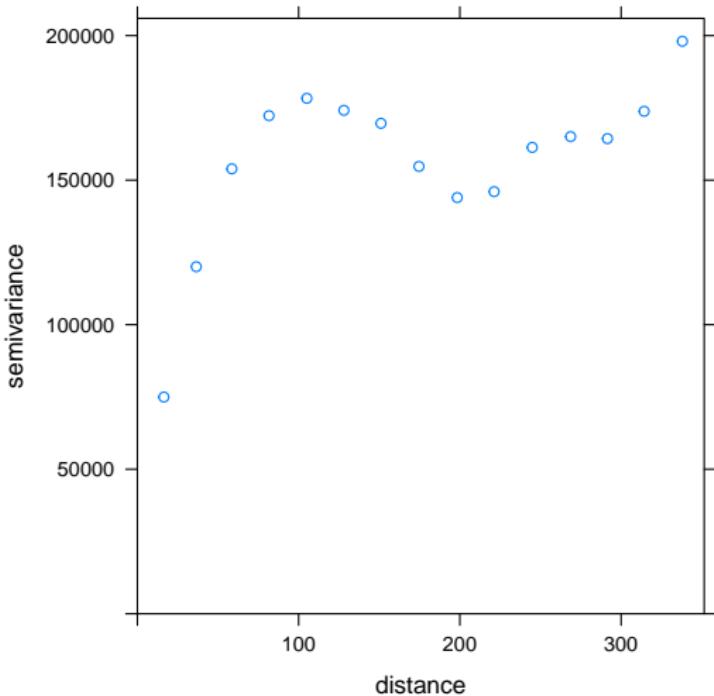


Figure 4 – Semi-variogram for the residuals of model.lm.

Example Las Rosas (21/23)

Moran correlogram : the value of the spatial lag at which I is no longer significantly positive can be used as an indication of the range of autocorrelation of the data.

- We create a list of neighbors using the k -nearest neighbors method
- We choose a row-standardised weights matrix.

```
library(spdep)
nlist <- knn2nb(knearneigh(Xutm, k=4))
I.d <- sp.correlogram(nlist, Xutm$resmodel.lm, order=10, method="I", style="W")
plot(I.d)
```

Example Las Rosas (22/23)

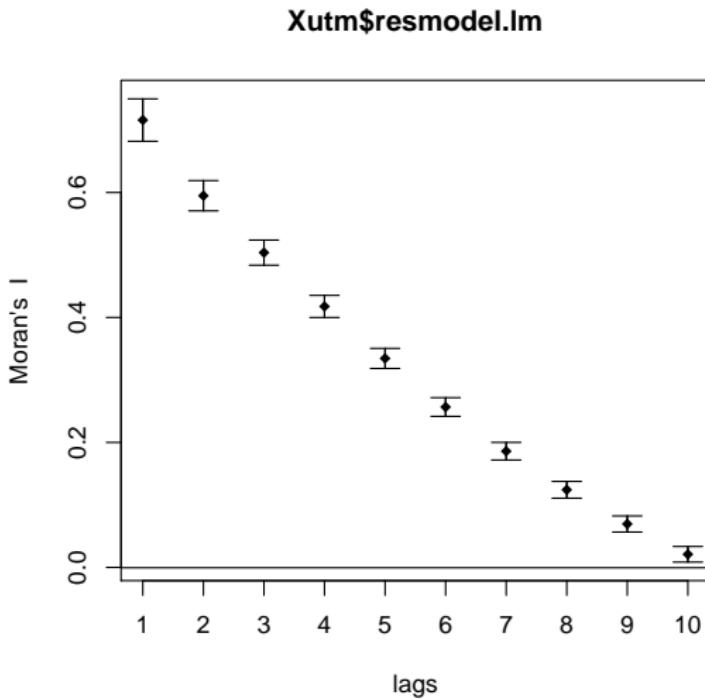


Figure 5 – Moran correlogram for residuals of model2.lm.

Example Las Rosas (23/23)

```
library(spdep)
nlist <- knn2nb(knearneigh(Xutm,k=8))
W <- nb2listw(nlist,style="W")
lm.morantest(model.lm,W)

##
## Global Moran I for regression residuals
##
## data:
## model: lm(formula = YIELD ~ accuf + N + slope + aspect + hshade +
## accuf * slope + accuf * aspect + accuf * hshade, data = Xutm@data)
## weights: W
##
## Moran I statistic standard deviate = 56.299, p-value < 2.2e-16
## alternative hypothesis: greater
## sample estimates:
## Observed Moran I      Expectation      Variance
##      0.6629755568     -0.0046606672     0.0001406322
```

Table of Contents

1 Origins and Consequences of Spatial Autocorrelation

2 Working example : Las Rosas

3 Spatial Lag Model

- Without explanatory variable
- With explanatory variable
- About the variance-covariance matrix of Y
- Fitting the model

4 Spatial Error Model

5 Choosing Between Spatial Lag, Error and SAC models

6 Extended Linear Models

7 Bibliography

Spatial lag model without explanatory variables

A spatial lag model with zero mean value and no explanatory variable has the form :

$$\begin{aligned} Y &= \rho WY + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I), \end{aligned} \tag{8}$$

where WY represents the spatial lag.

Interpretation

The value of Y at one location is directly associated with the values of the process Y at nearby locations.

For instance high productivity of a plant at one location is associated with high productivity at nearby locations (but there is no notion of causality).

Spatial lag model with explanatory variables

A spatial lag model with explanatory variables :

$$\begin{aligned} Y &= \rho WY + X\beta + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I). \end{aligned} \tag{9}$$

Interpretation

This model can be interpreted using different points of view (Anselin 1992).

- 1** Specification of the spatial weights matrix W and estimation of ρ are indicators of the **nature and strength of spatial interaction**.
- 2** $Y = (I - \rho W)^{-1}(X\beta + \epsilon)$, and $\mathbb{E}(Y) = (I - \rho W)^{-1}X\beta$: non-linear effect of the spatial autocorrelation on the expected value of Y . **The influence of the spatial structure is modelled through the error term and through the explanatory variables (influence of the neighborhood).**

Prediction $\hat{Y} = (I - \hat{\rho}W)^{-1}\hat{X}\hat{\beta}$ is mainly driven by the neighborhood. If we use $\hat{Y} = \hat{X}\hat{\beta}$, we have a bias $-(\rho W)^{-1}X\beta$.

About the variance-covariance matrix of Y

$$\text{var}(Y) = \sigma^2(I - \rho W)^{-1}(I - \rho W')^{-1}.$$

- Impacted by the magnitude of the variance of the error term σ^2 .
- Impacted by the spatial structure through the term $(I - \rho W)^{-1}(I - \rho W')^{-1}$.
- Enforced by the model, we do not have to specify it. **The spatial autocorrelation structure of Y is enforced by the model.**

Fitting the model

Estimation of the parameters β , σ^2 and ρ

Maximum likelihood approach.

In practice

The expressions of $\hat{\beta}$, $\hat{\sigma}^2$ and $\hat{\rho}$ that maximise the likelihood are not easy to obtain (it would be much easier if ρ was known)

⇒ use of a **numerical scheme** analogous to the Newton-Raphson method :

- A value of $\hat{\rho}$ is fixed.
- Maximum likelihood estimates $\hat{\beta}$ and $\hat{\sigma}^2$ calculated with $\hat{\rho}$ fixed.
- The two preceding steps are iterated : another value of $\hat{\rho}$ increasing the likelihood is fixed, $\hat{\beta}$ and $\hat{\sigma}^2$ are calculated to maximise the likelihood, then fix $\hat{\rho}$ again,...

Spatial lag model : Las Rosas (1/4)

```

nlist <- knn2nb(knearneigh(Xutm, k=8))
W <- nb2listw(nlist, style="W")
summary(Xutm@data[, 1:8])

##          YIELD             N            elev        slope
##  Min.   : 574.2   Min.   : 0.00   Min.   :268.0   Min.   :0.007235
##  1st Qu.:2290.1   1st Qu.: 39.00   1st Qu.:271.0   1st Qu.:0.016802
##  Median :3826.5   Median : 50.60   Median :273.6   Median :0.024580
##  Mean   :3412.5   Mean   : 64.93   Mean   :273.7   Mean   :0.024072
##  3rd Qu.:4511.8   3rd Qu.: 99.80   3rd Qu.:276.2   3rd Qu.:0.031336
##  Max.   :5347.9   Max.   :124.60   Max.   :280.0   Max.   :0.041564
##          slopeX           accu         aspect      hshade
##  Min.   :-23.7891   Min.   :-143.137   Min.   :1.850   Min.   :0.8509
##  1st Qu.:-11.8105   1st Qu.:-67.273   1st Qu.:2.108   1st Qu.:0.8589
##  Median : -1.7972   Median : 20.225   Median :2.964   Median :0.8625
##  Mean   :  0.2142   Mean   : -4.171   Mean   :3.240   Mean   :0.8613
##  3rd Qu.: 13.4638   3rd Qu.: 56.062   3rd Qu.:4.428   3rd Qu.:0.8637
##  Max.   : 21.8369   Max.   : 89.600   Max.   :4.792   Max.   :0.8669

```

Spatial lag model : Las Rosas (2/4)

```
Xutm$YIELD_scaled <- (Xutm$YIELD-mean(Xutm$YIELD))/sd(Xutm$YIELD)
# Xutm$N_scaled <- (Xutm$N-mean(Xutm$N))/sd(Xutm$N)
Xutm$slope_scaled <- (Xutm$slope-mean(Xutm$slope))/sd(Xutm$slope)
# Xutm$aspect_scaled <- (Xutm$aspect-mean(Xutm$aspect))/sd(Xutm$aspect)
Xutm$hshade_scaled <- (Xutm$hshade-mean(Xutm$hshade))/sd(Xutm$hshade)

library(spatialreg)
f <- as.formula("YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_
                  + accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled")
mod.lag <- lagsarlm(f, data=Xutm, listw=W)
```

Interpretation of the spatial lag model : there is some interaction (competition) between corn plants.

Spatial lag model : Las Rosas (3/4)

```
summary(mod.lag)
```

```
##  
## Call:lagsarlm(formula = f, data = Xutm, listw = W)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -1.3179579 -0.1180341  0.0029359  0.1081311  0.6812413  
##  
## Type: lag  
## Coefficients: (asymptotic standard errors)  
##  
##             Estimate Std. Error z value Pr(>|z|)  
## (Intercept) 0.02004197 0.02669322 0.7508 0.4527573  
## accuflow    -0.20148029 0.05537151 -3.6387 0.0002740  
## N           0.00165631 0.00011166 14.8340 < 2.2e-16  
## slope_scaled -0.05422472 0.00865827 -6.2628 3.782e-10  
## aspect      -0.01605523 0.00779494 -2.0597 0.0394273  
## hshade_scaled -0.03632647 0.01541618 -2.3564 0.0184537
```

Spatial lag model : Las Rosas (4/4)

```
## accuflow:slope_scaled  0.06974847  0.01518287  4.5939  4.351e-06
## accuflow:aspect        -0.00077238  0.01448327 -0.0533  0.9574698
## accuflow:hshade_scaled 0.06926005  0.01876857  3.6902  0.0002241
##
## Rho: 0.8753, LR test value: 1919.8, p-value: < 2.22e-16
## Asymptotic standard error: 0.01351
##      z-value: 64.788, p-value: < 2.22e-16
## Wald statistic: 4197.5, p-value: < 2.22e-16
##
## Log likelihood: 346.279 for lag model
## ML residual variance (sigma squared): 0.033335, (sigma: 0.18258)
## Number of observations: 1704
## Number of parameters estimated: 11
## AIC: -670.56, (AIC for lm: 1247.2)
## LM test for residual autocorrelation
## test value: 32.674, p-value: 1.0897e-08
```

Table of Contents

1 Origins and Consequences of Spatial Autocorrelation

2 Working example : Las Rosas

3 Spatial Lag Model

4 Spatial Error Model

- Formulation
- About the variance-covariance matrix of Y and fitting

5 Choosing Between Spatial Lag, Error and SAC models

6 Extended Linear Models

7 Bibliography

Spatial Error Model (1/2)

$$\begin{aligned} Y &= X\beta + \eta && (10) \\ \eta &= \lambda W\eta + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I). \end{aligned}$$

Interpretation

- 1 Like a classical linear model, but with a correlated structure for the error term.

This autocorrelation is generally considered to be a nuisance : the primary interest is often the relationship between the explanatory variables X and the response variable Y .

The spatial autocorrelation is just taken into account through the error term.

Spatial Error Model (2/2)

- 2 The influence of the spatial structure is modelled only on the error term $Y = X\beta + (I - \lambda W)^{-1}\epsilon$.

Prediction $\hat{Y} = X\hat{\beta}$ is driven by the values of the explanatory variables at the location for which we want the prediction. Be careful, to have an unbiased estimation of β , you must use the spatial error model and not the classical linear model if your data are driven by this spatial error model.

This model can be written as a classical linear model :

$$\begin{aligned} Y - \lambda WY &= X\beta - \lambda WY + \eta \\ &= (X - \lambda WX)\beta + \epsilon \\ \tilde{Y} &= \tilde{X}\beta + \epsilon \end{aligned}$$

About the variance-covariance matrix of Y and fitting

$$\text{var}(Y) = \sigma^2(I - \lambda W)^{-1}(I - \lambda W')^{-1}. \quad (11)$$

- Impacted by the magnitude of the variance of the error term σ^2 .
- Impacted by the spatial structure through the term $(I - \lambda W)^{-1}(I - \lambda W')^{-1}$.
- Enforced by the model, we do not have to specify it. **The spatial autocorrelation structure of Y is enforced by the model.**

Fitting the model

The approach is the same as for the spatial lag model, with ρ replaced by λ .

Spatial error model : Las Rosas (1/2)

Interpretation of the spatial error model : the yield at one location is supposed mainly driven by the values of the explanatory variables at this location, these are the error terms which are spatially autocorrelated.

```
mod.err <- errorsarlm(f,data=Xutm,listw=W)
summary(mod.err)
```

```
## Type: error
## Coefficients: (asymptotic standard errors)
##                                     Estimate Std. Error z value Pr(>|z|)
## (Intercept)                 1.13645551  0.28442272  3.9957 6.452e-05
## accuflow                  -1.34933604  0.43064505 -3.1333 0.0017286
## N                          0.00196741  0.00011066 17.7797 < 2.2e-16
## slope_scaled                -0.42523669  0.06866555 -6.1929 5.908e-10
## aspect                      -0.33772997  0.08921225 -3.7857 0.0001533
## hshade_scaled                0.36372642  0.12239612  2.9717 0.0029614
```

Spatial error model : Las Rosas (2/2)

```
## accuflow:slope_scaled    0.58286320  0.09749498  5.9784 2.254e-09
## accuflow:aspect           0.15722109  0.13008315  1.2086 0.2268089
## accuflow:hshade_scaled   -0.09240629  0.14049648 -0.6577 0.5107229
##
## Lambda: 0.93964, LR test value: 1934.2, p-value: < 2.22e-16
## Asymptotic standard error: 0.009443
##      z-value: 99.506, p-value: < 2.22e-16
## Wald statistic: 9901.5, p-value: < 2.22e-16
##
## Log likelihood: 353.489 for error model
## ML residual variance (sigma squared): 0.031593, (sigma: 0.17774)
## Number of observations: 1704
## Number of parameters estimated: 11
## AIC: -684.98, (AIC for lm: 1247.2)
```

Table of Contents

- 1 Origins and Consequences of Spatial Autocorrelation
- 2 Working example : Las Rosas
- 3 Spatial Lag Model
- 4 Spatial Error Model
- 5 Choosing Between Spatial Lag, Error and SAC models
- 6 Extended Linear Models
- 7 Bibliography

Choosing Between Spatial Lag, Error and SAC models (1/2)

Spatial autocorrelation detected in residuals of a classical linear model

- ⇒ take into account this autocorrelation
- ⇒ choose between spatial lag model and spatial error model (or an extended linear model).

These two models can be combined in a SAC/SARAR model :

$$\begin{aligned} Y &= \rho W_1 Y + X\beta + \eta & (12) \\ \eta &= \lambda W_2 \eta + \epsilon \\ \epsilon &\sim \mathcal{N}(0, \sigma^2 I), \end{aligned}$$

where W_1 can be equal to W_2 and X cannot simply be a vector of ones (for identifiability).

Choosing Between Spatial Lag, Error and SAC models (2/2)

$$\text{1st test : } \begin{cases} H_0 : & \rho = 0, \\ H_1 : & \rho \neq 0 \end{cases} \quad \text{and} \quad \text{2nd test : } \begin{cases} H_0 : & \lambda = 0, \\ H_1 : & \lambda \neq 0 \end{cases}$$

- H_0 kept for both tests \Rightarrow keep a **classical linear model**, there is no spatial autocorrelation of the residuals.
- H_0 kept for the second test and H_1 non-rejected for the first test \Rightarrow **spatial lag model**.
- H_0 kept for the first test and H_1 non-rejected for the second test \Rightarrow **spatial error model**.
- H_1 non-rejected for one of these tests \Rightarrow **try to fit a SAC model** and compare it to a spatial lag or spatial error model.

In practice :

- These tests are carried using the maximum likelihood approach.
- Another possibility : use the AIC criteria.

Example Las Rosas (1/31)

```
f <- as.formula("YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled")  
  
model.lm_scaled <- lm(f, data=Xutm)
mod.lag <- lagsarlm(f, data=Xutm, listw=W)
mod.err <- errorsarlm(f, data=Xutm, listw=W)
LR.sarlm(model.lm_scaled, mod.lag)  
  
##  
## Likelihood ratio for spatial linear models  
##  
## data:  
## Likelihood ratio = -1919.8, df = 1, p-value < 2.2e-16  
## sample estimates:  
## Log likelihood of model.lm_scaled           Log likelihood of mod.lag  
##                                         -613.6159          346.2790
```

Example Las Rosas (2/31)

```
LR.sarlm(model.lm_scaled,mod.err)

##
## Likelihood ratio for spatial linear models
##
## data:
## Likelihood ratio = -1934.2, df = 1, p-value < 2.2e-16
## sample estimates:
## Log likelihood of model.lm_scaled           Log likelihood of mod.err
##                               -613.6159                  353.4890

AIC(mod.lag,mod.err)

##          df      AIC
## mod.lag 11 -670.5580
## mod.err 11 -684.9779
```

Example Las Rosas (3/31)

AIC \Rightarrow we prefer the spatial error model.

As the H_1 hypothesis is not rejected for both tests, we can fit a SAC model.

```
mod.sac <- sacsarlm(f, data=Xutm, listw=W)
LR.sarlm(mod.sac, mod.lag)

##
## Likelihood ratio for spatial linear models
##
## data:
## Likelihood ratio = 104.3, df = 1, p-value < 2.2e-16
## sample estimates:
## Log likelihood of mod.sac Log likelihood of mod.lag
##                  398.4283                 346.2790
```

Example Las Rosas (4/31)

```
LR.sarlm(mod.sac,mod.err)
```

```
##  
## Likelihood ratio for spatial linear models  
##  
## data:  
## Likelihood ratio = 89.879, df = 1, p-value < 2.2e-16  
## sample estimates:  
## Log likelihood of mod.sac Log likelihood of mod.err  
## 398.4283 353.4890
```

```
AIC(mod.lag,mod.err,mod.sac)
```

```
## df AIC  
## mod.lag 11 -670.5580  
## mod.err 11 -684.9779  
## mod.sac 12 -772.8566
```

Example Las Rosas (5/31)

The SAC model is better than the spatial lag or the spatial error model : we both have competition between corn plants, and the error terms are spatially correlated.

Improve the spatial SAC model by performing model selection.

- Try to remove explanatory variables or interactions between them and to include variables which are not present in `mod.err`.
- Likelihood ratio test or AIC criteria to select the best model.

Example Las Rosas (6/31)

```
summary(mod.sac)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)           1.75765004  0.77749615  2.2607  0.023781
## accuflow            -1.16490864  0.61228927 -1.9025  0.057100
## N                   0.00151935  0.00011262 13.4915 < 2.2e-16
## slope_scaled        -0.17491748  0.13328369 -1.3124  0.189395
## aspect              -0.52703545  0.20278078 -2.5990  0.009348
## hshade_scaled       0.38036319  0.21454549  1.7729  0.076249
## accuflow:slope_scaled 0.18923145  0.14064994  1.3454  0.178494
## accuflow:aspect      0.26730066  0.17766099  1.5046  0.132439
## accuflow:hshade_scaled -0.26141594  0.18786214 -1.3915  0.164065
```

Example Las Rosas (7/31)

```
f2 <- as.formula("YIELD_scaled ~ accuf + N + slope_scaled + aspect  
                  + hshade_scaled + accuf*aspect + accuf*hshade_scaled")  
mod.sac2 <- sacsarlm(f2,data=Xutm,listw=W)
```

```
LR.sarlm(mod.sac, mod.sac2)
```

```
##  
## Likelihood ratio for spatial linear models  
##  
## data:  
## Likelihood ratio = 1.6409, df = 1, p-value = 0.2002  
## sample estimates:  
## Log likelihood of mod.sac Log likelihood of mod.sac2  
##                      398.4283                      397.6078
```

Example Las Rosas (8/31)

```
summary(mod.sac2)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                                     Estimate Std. Error z value Pr(>|z|)
## (Intercept)           1.66913033  0.81695206  2.0431  0.04104
## accuflow            -0.99822066  0.60807547 -1.6416  0.10067
## N                   0.00151476  0.00011243 13.4726 < 2e-16
## slope_scaled        -0.10539026  0.12812374 -0.8226  0.41075
## aspect              -0.50983334  0.20672089 -2.4663  0.01365
## hshade_scaled       0.36139428  0.21795862  1.6581  0.09730
## accuflow:aspect     0.28738455  0.17766510  1.6176  0.10576
## accuflow:hshade_scaled -0.30201642  0.18562952 -1.6270  0.10374
```

Example Las Rosas (9/31)

```
f3 <- as.formula("YIELD_scaled ~ accuf + N + slope_scaled + aspect  
+ hshade_scaled + accuf*hshade_scaled")  
mod.sac3 <- sacsarlm(f3,data=Xutm,listw=W)
```

```
LR.sarlm(mod.sac2, mod.sac3)
```

```
##  
## Likelihood ratio for spatial linear models  
##  
## data:  
## Likelihood ratio = 2.572, df = 1, p-value = 0.1088  
## sample estimates:  
## Log likelihood of mod.sac2 Log likelihood of mod.sac3  
## 397.6078 396.3219
```

Example Las Rosas (10/31)

```
summary(mod.sac3)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                                     Estimate Std. Error z value Pr(>|z|)
## (Intercept)                 1.1470900  0.7703550  1.4890  0.13648
## accuflow                  -0.0170678  0.0496118 -0.3440  0.73083
## N                          0.0015136  0.0001124 13.4661 < 2e-16
## slope_scaled                -0.0824190  0.1285267 -0.6413  0.52135
## aspect                      -0.3345291  0.1766718 -1.8935  0.05829
## hshade_scaled                0.1320709  0.1711735  0.7716  0.44037
## accuflow:hshade_scaled   -0.0085828  0.0390317 -0.2199  0.82595
```

Example Las Rosas (11/31)

```
f4 <- as.formula("YIELD_scaled ~ accuf + N + slope_scaled + aspect  
+ hshade_scaled")  
mod.sac4 <- sacsarlm(f4,data=Xutm,listw=W)
```

```
LR.sarlm(mod.sac3, mod.sac4)
```

```
##  
## Likelihood ratio for spatial linear models  
##  
## data:  
## Likelihood ratio = 0.04834, df = 1, p-value = 0.826  
## sample estimates:  
## Log likelihood of mod.sac3 Log likelihood of mod.sac4  
## 396.3219 396.2977
```

Example Las Rosas (12/31)

```
summary(mod.sac4)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                  Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.14654210  0.76990060  1.4892  0.13643
## accuflow      -0.01694282  0.04960867 -0.3415  0.73271
## N              0.00151400  0.00011239 13.4715 < 2e-16
## slope_scaled   -0.08138525  0.12841259 -0.6338  0.52622
## aspect         -0.33388599  0.17662274 -1.8904  0.05871
## hshade_scaled   0.12672191  0.16932172  0.7484  0.45421
```

Example Las Rosas (13/31)

```
f5 <- as.formula("YIELD_scaled ~ N + slope_scaled + aspect + hshade_scaled")
mod.sac5 <- sacsarlm(f5, data=Xutm, listw=W)

LR.sarlm(mod.sac4, mod.sac5)

##
## Likelihood ratio for spatial linear models
##
## data:
## Likelihood ratio = 0.11589, df = 1, p-value = 0.7335
## sample estimates:
## Log likelihood of mod.sac4 Log likelihood of mod.sac5
##                      396.2977                  396.2397
```

Example Las Rosas (14/31)

```
summary(mod.sac5)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.14286973  0.77325648  1.4780  0.13941
## N                 0.00151338  0.00011237 13.4679 < 2e-16
## slope_scaled   -0.08042127  0.12851282 -0.6258  0.53146
## aspect          -0.33420918  0.17686740 -1.8896  0.05881
## hshade_scaled   0.12711273  0.16947486  0.7500  0.45323
```

Example Las Rosas (15/31)

```
f6 <- as.formula("YIELD_scaled ~ N + aspect + hshade_scaled")
mod.sac6 <- sacsarlm(f6, data=Xutm, listw=W)
```

```
LR.sarlm(mod.sac5, mod.sac6)
```

```
##
## Likelihood ratio for spatial linear models
##
## data:
## Likelihood ratio = 0.37328, df = 1, p-value = 0.5412
## sample estimates:
## Log likelihood of mod.sac5 Log likelihood of mod.sac6
##                      396.2397                  396.0531
```

Example Las Rosas (16/31)

```
summary(mod.sac6)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                   Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.26179451  0.77180194  1.6349  0.10208
## N              0.00151327  0.00011229 13.4770 < 2e-16
## aspect        -0.36893656  0.16958743 -2.1755  0.02959
## hshade_scaled  0.14790879  0.16488224  0.8971  0.36969
```

Example Las Rosas (17/31)

```
f7 <- as.formula("YIELD_scaled ~ N + aspect")
mod.sac7 <- sacsarlm(f7,data=Xutm,listw=W)

LR.sarlm(mod.sac6, mod.sac7)

##
## Likelihood ratio for spatial linear models
##
## data:
## Likelihood ratio = 0.73526, df = 1, p-value = 0.3912
## sample estimates:
## Log likelihood of mod.sac6 Log likelihood of mod.sac7
##                      396.0531                  395.6855
```

Example Las Rosas (18/31)

```
summary(mod.sac7)
```

```
## Type: sac
## Coefficients: (asymptotic standard errors)
##                  Estimate Std. Error z value Pr(>|z|)
## (Intercept)  1.23727293  0.80598355  1.5351  0.12476
## N            0.00150883  0.00011217 13.4515 < 2e-16
## aspect      -0.35675187  0.17054142 -2.0919  0.03645
```

Example Las Rosas (19/31)

We choose the model mod.sac7. The same conclusion is obtained using the AIC.

```
AIC(mod.sac,mod.sac2,mod.sac3,mod.sac4,mod.sac5,mod.sac6,mod.sac7)
```

```
##           df      AIC
## mod.sac   12 -772.8566
## mod.sac2  11 -773.2157
## mod.sac3  10 -772.6437
## mod.sac4   9 -774.5954
## mod.sac5   8 -776.4795
## mod.sac6   7 -778.1062
## mod.sac7   6 -779.3710
```

This model is quite simple compared to the classical linear model obtained before!

We have to check that the assumptions are verified on the residuals of mod.sac7.

Example Las Rosas (20/31)

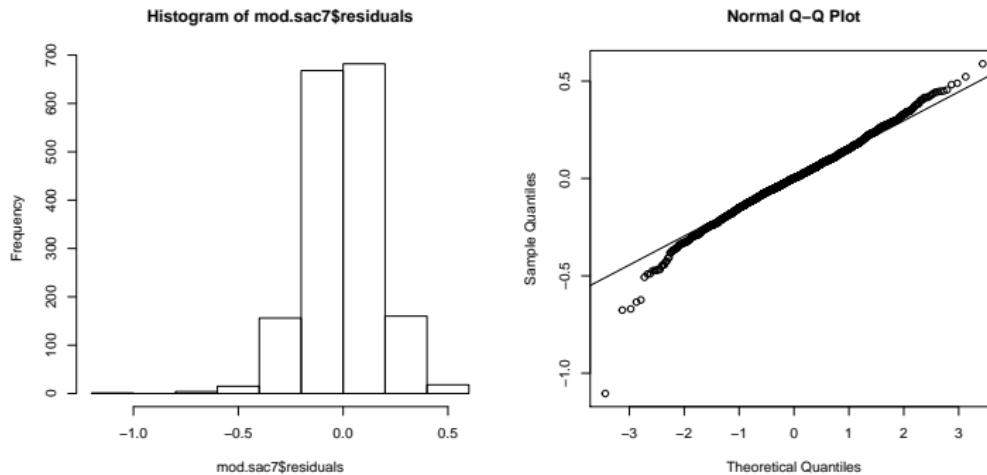


Figure 6 – Histogram of residuals of mod.sac7 and associated QQ plot.

Example Las Rosas (21/31)

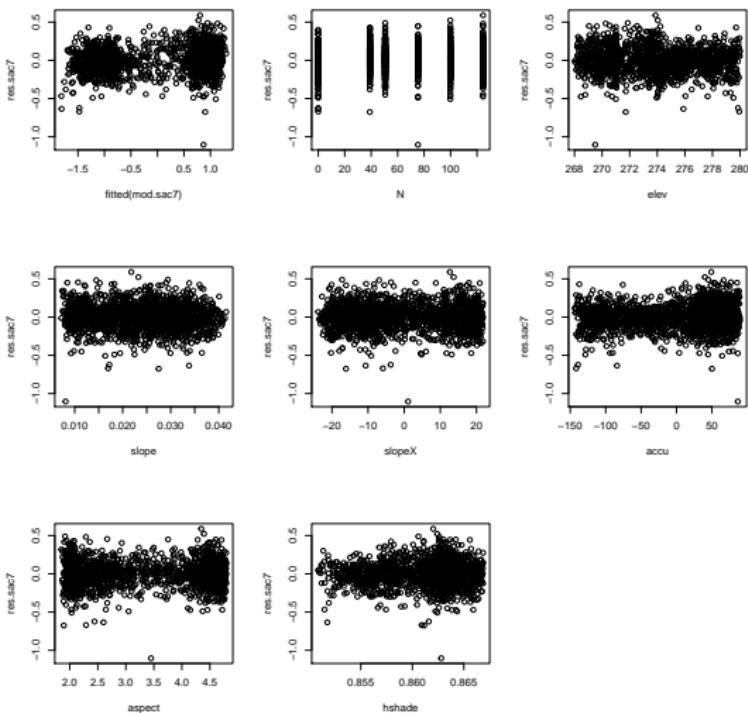


Figure 7 – Residuals of mod.sac7 against every possible explanatory variable.

Example Las Rosas (22/31)

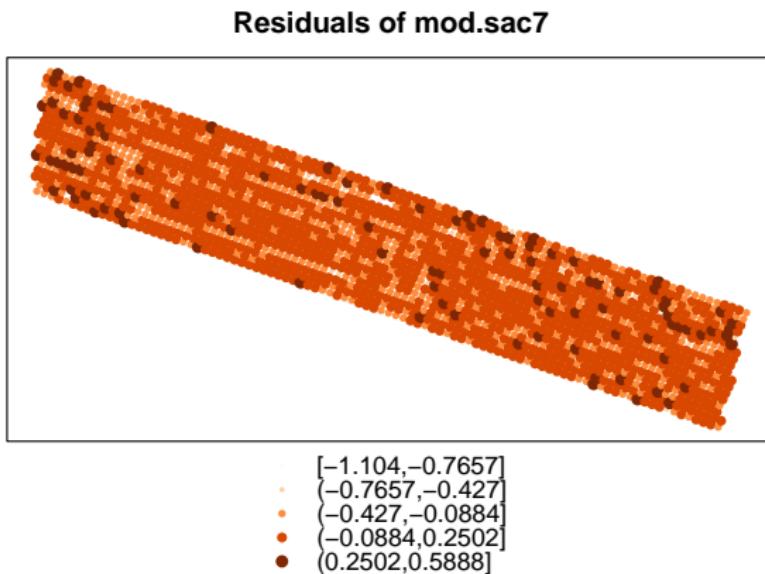


Figure 8 – Bubble map for residuals of mod.sac7.

Example Las Rosas (23/31)

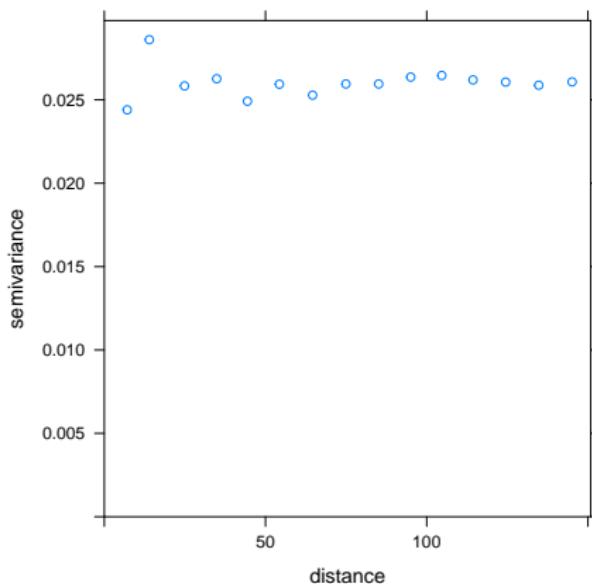


Figure 9 – Semi-variogram for the residuals of mod.sac7.

Example Las Rosas (24/31)

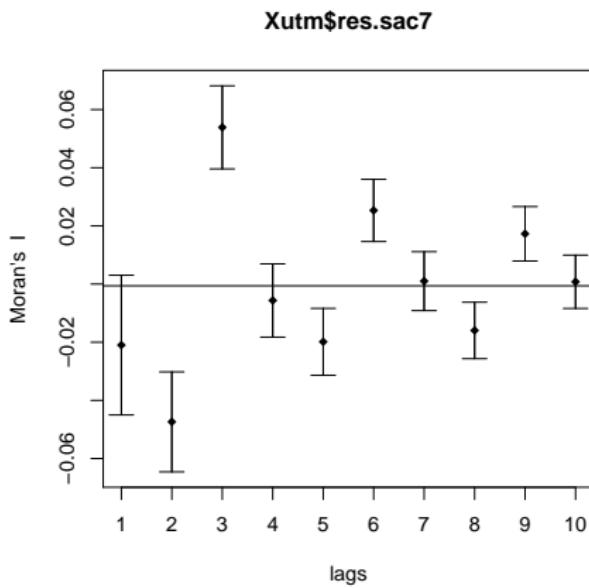


Figure 10 – Moran correlogram for residuals of mod.sac7.

Example Las Rosas (25/31)

```
moran.mc(Xutm@data$res.sac7,W,nsim=1000,alternative="greater")  
  
##  
## Monte-Carlo simulation of Moran I  
##  
## data: Xutm@data$res.sac7  
## weights: W  
## number of simulations + 1: 1001  
##  
## statistic = -0.020985, observed rank = 34, p-value = 0.966  
## alternative hypothesis: greater
```

Example Las Rosas (26/31)

```
moran.mc(Xutm@data$res.sac7,W,nsim=1000,alternative="less")  
  
##  
## Monte-Carlo simulation of Moran I  
##  
## data: Xutm@data$res.sac7  
## weights: W  
## number of simulations + 1: 1001  
##  
## statistic = -0.020985, observed rank = 36, p-value = 0.03596  
## alternative hypothesis: less
```

Example Las Rosas (27/31)

Predictions using this spatial SAC model can be done.

Predictions for the data on which the model was fitted (fitted values) :

```
pred <- as.data.frame(predict.sarlm(mod.sac7, listw=W, pred.type="trend"))
head(pred)

##           fit
## 1 -0.1661076
## 2 -0.1783937
## 3 -0.1903548
## 4 -0.1979362
## 5 -0.2069975
## 6 -0.2133545
```

Example Las Rosas (28/31)

Previsions for new data : we want to predict the yield in the field if the nitrogen content is increased uniformly by one unit (using fertilizer).

For spatial-lag model :

```
Xutm2 <- Xutm  
Xutm2@data$N <- Xutm@data$N + 1  
newpred <- as.data.frame(predict.sarlm(mod.lag, newdata=Xutm2@data,  
                                listw = W, pred.type="TS"))  
head(newpred)  
  
##          fit      trend    signal  
## 1 0.4821111 0.14613577 0.3359753  
## 2 0.4546835 0.13409381 0.3205896  
## 3 0.4286981 0.12223805 0.3064601  
## 4 0.3613864 0.11290496 0.2484815  
## 5 0.3086342 0.10375643 0.2048778  
## 6 0.2627418 0.09649667 0.1662452
```

Example Las Rosas (29/31)

For spatial-error model :

```
newpred <- as.data.frame(predict.sarlm(mod.err, newdata=utm2@data,
                                         listw = W, pred.type="TS"))

head(newpred)

##           fit      trend signal
## 1  0.20902349  0.20902349     0
## 2  0.11276585  0.11276585     0
## 3  0.01799241  0.01799241     0
## 4 -0.05881832 -0.05881832     0
## 5 -0.11794635 -0.11794635     0
## 6 -0.13481881 -0.13481881     0
```

Example Las Rosas (30/31)

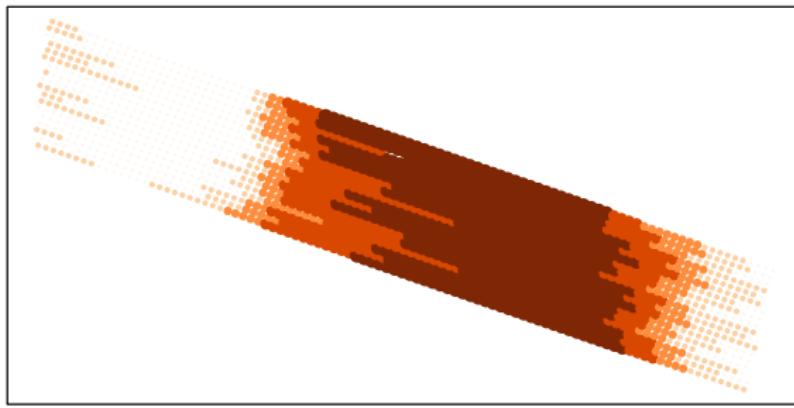
For spatial SAC model :

```
newpred <- as.data.frame(predict.sarlm(mod.sac7, newdata=Xutm2@data,
                                         listw = W, pred.type="trend"))
head(newpred)

##          fit
## 1 -0.1645987
## 2 -0.1768848
## 3 -0.1888460
## 4 -0.1964274
## 5 -0.2054887
## 6 -0.2118457
```

Example Las Rosas (31/31)

**Predicted yield with an increase
of nitrogen content of one, using the spatial SAC model**



- $[-0.4678, -0.2259]$
- $(-0.2259, 0.01597]$
- $(0.01597, 0.2579]$
- $(0.2579, 0.4998]$
- $(0.4998, 0.7417]$

Figure 11 – Bubble map for the predicted yield when N is increased by one unit, for the spatial SAC model.

Table of Contents

1 Origins and Consequences of Spatial Autocorrelation

2 Working example : Las Rosas

3 Spatial Lag Model

4 Spatial Error Model

5 Choosing Between Spatial Lag, Error and SAC models

6 Extended Linear Models

- Classical Linear Model versus Extended Linear Model
- Modelling Spatial Correlation

7 Bibliography

Classical linear model versus extended linear model (1/2)

Y quantitative variable to explain, explanatory variables quantitative or qualitative :

$$Y = X\beta + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I). \quad (13)$$

Possible extensions

On the variance-covariance matrix of the residuals (among others).

In classical linear models :

- The residuals (therefore the observations) are supposed independent.
- The residuals are supposed homoscedastic.

Classical linear model versus extended linear model (2/2)

$$Y = X\beta + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \Lambda). \quad (14)$$

- 1 If Λ is diagonal, but with varying coefficients on the diagonal, \Rightarrow heteroscedasticity.
- 2 If Λ has non-null coefficients outside the diagonal \Rightarrow correlation between the residuals, dependence structure of the residuals. This dependence can be temporal, spatial or more general.

In practice, once a modelisation has been chosen

- Parameters of these extended linear models (regression coefficients and coefficients of the variance-covariance matrix) estimated using the maximum likelihood estimators.
- Numerically obtained by solving an ordinary least-squares problem.

Modelling Spatial Correlation

Extended linear model vs regression models for spatially autocorrelated data

- Models for spatially autocorrelated data : special cases of extended linear models.
- In extended linear models Λ can take any form.
- In the regression models designed for spatially autocorrelated data, the form of Λ is enforced by the model.
- Regression models for spatially autocorrelated data often more intuitive than extended linear models.

Choosing the modelisation of the spatial dependency (1/2)

Look at the form of the semi-variogram

- Choosing the form of the variance-covariance matrix $\Lambda \Leftrightarrow$ to choose a semi-variogram pattern.
- The form of the empirical semi-variogram can guide us to choose a semi-variogram pattern.

Use classical model selection methods

AIC, BIC, tests between nested models.

Choosing the modelisation of the spatial dependency (2/2)

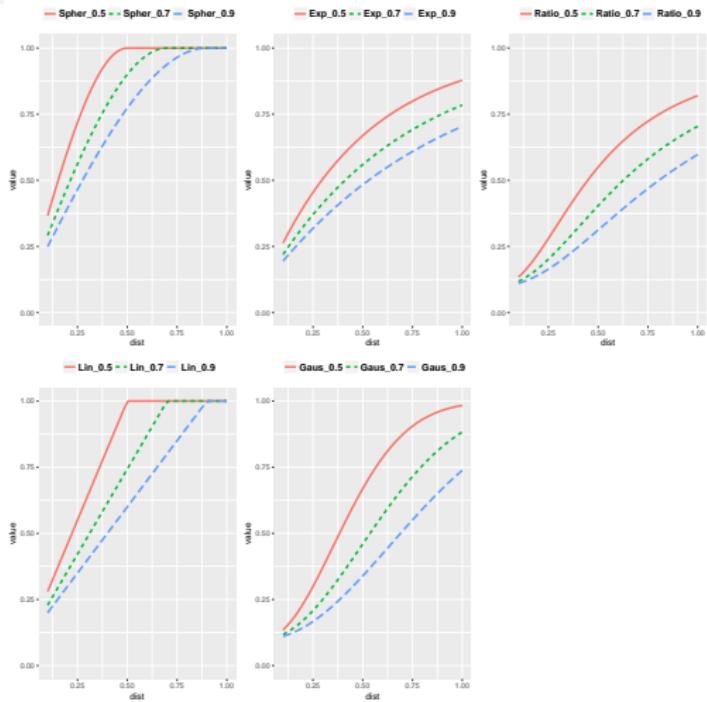


Figure 12 – Different semi-variogram patterns : Spherical, Exponential, Rational quadratic, Linear and Gaussian. Each pattern has a nugget of 0.1. The value of

Example Las Rosas (1/8)

```
library(nlme)
model2.lm <- gls(YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled,
modSpher <- gls(YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_s
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled,
data=Xutm, correlation=corSpher(form=~x+y,nugget=T))
modLin <- gls(YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_sca
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled,
data=Xutm, correlation=corLin(form=~x+y,nugget=T))
modRatio <- gls(YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_s
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled,
data=Xutm, correlation=corRatio(form=~x+y,nugget=T))
modGaus <- gls(YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_s
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled,
data=Xutm, correlation=corGaus(form=~x+y,nugget=T))
modExp <- gls(YIELD_scaled ~ accuf + N + slope_scaled + aspect + hshade_sca
+ accuf*slope_scaled + accuf*aspect + accuf*hshade_scaled,
data=Xutm, correlation=corExp(form=~x+y,nugget=T))
```

Example Las Rosas (2/8)

```
AIC(modSpher,modLin,modRatio,modGaus,modExp)
```

```
##           df      AIC
## modSpher 12 -837.72045
## modLin   12  -14.27035
## modRatio 12 -785.46606
## modGaus  12 -752.27858
## modExp   12 -836.82643
```

```
VarioSpher_raw <- Variogram(modSpher, form =~ LONGITUDE + LATITUDE,
                           robust = TRUE, maxDist = 350, resType = "pearson")
VarioSpher_normalized <- Variogram(modSpher, form =~ LONGITUDE + LATITUDE,
                           robust = TRUE, maxDist = 350, resType = "normalized")
```

Example Las Rosas (3/8)

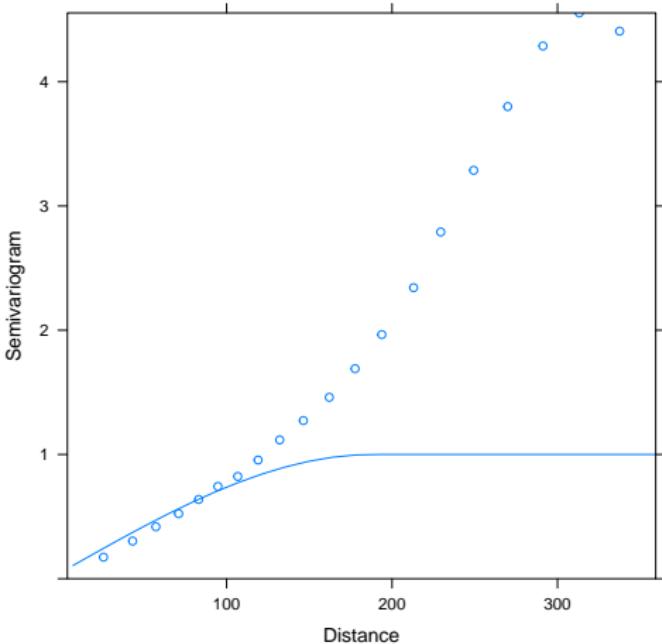


Figure 13 – Semi-variogram for the raw residuals of modSpher.

Example Las Rosas (4/8)

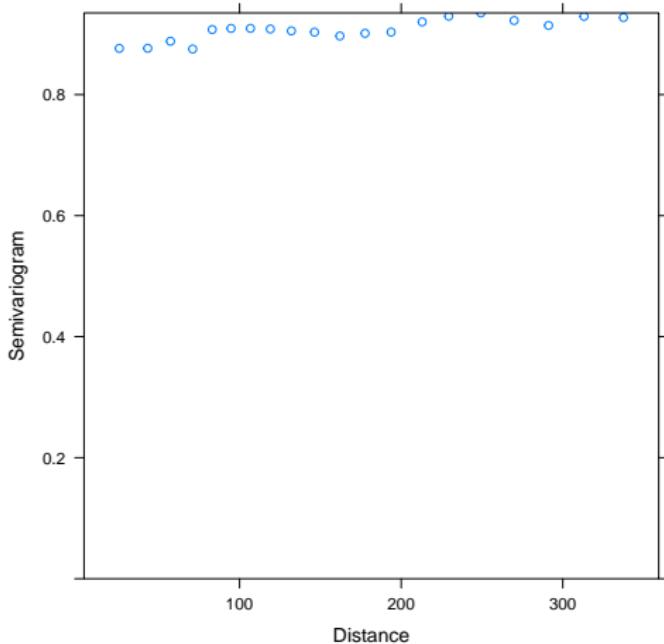
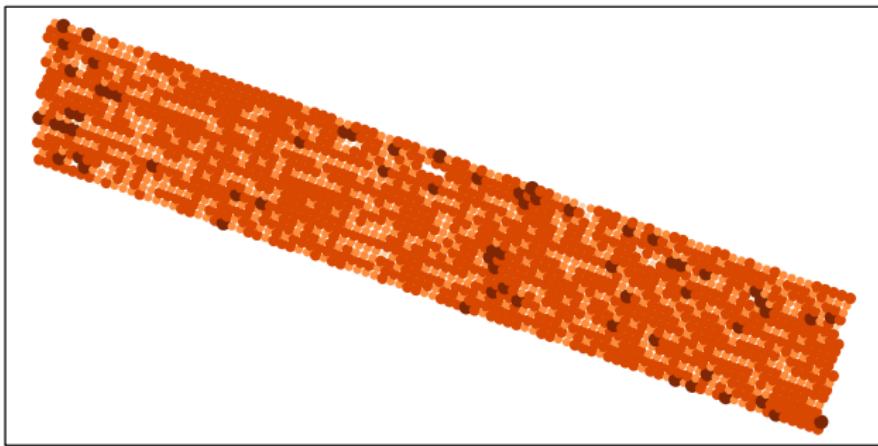


Figure 14 – Semi-variogram for the studentized residuals of modSpher.

Example Las Rosas(5/8)

Normalized residuals of modSpher

- $[-7.022, -4.812]$
- $(-4.812, -2.603]$
- $(-2.603, -0.3931]$
- $(-0.3931, 1.816]$
- $(1.816, 4.026]$

Figure 15 – Bubble map for residuals of modSpher.

Example Las Rosas (6/8)

```
moran.mc(Xutm@data$resSpherNorm,W,nsim=1000,alternative="greater")  
  
##  
## Monte-Carlo simulation of Moran I  
##  
## data: Xutm@data$resSpherNorm  
## weights: W  
## number of simulations + 1: 1001  
##  
## statistic = -0.023487, observed rank = 30, p-value = 0.97  
## alternative hypothesis: greater
```

Example Las Rosas (7/8)

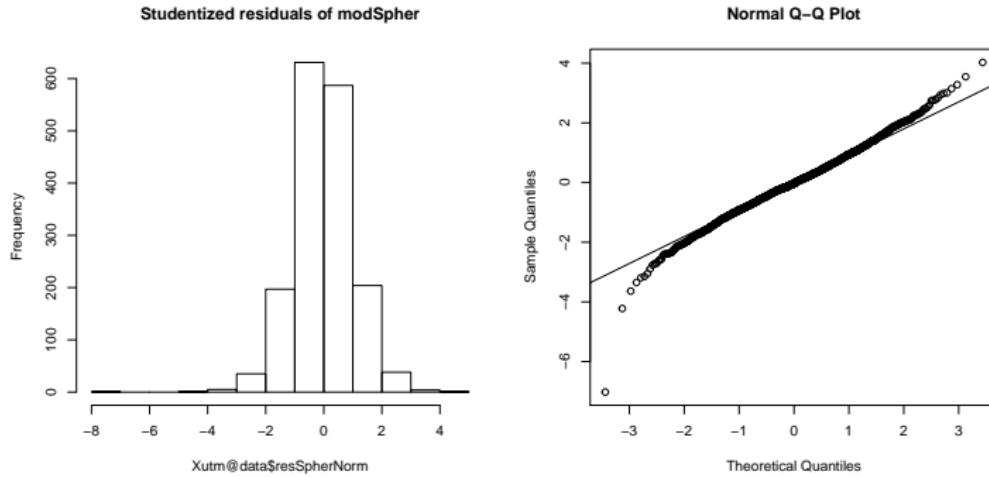


Figure 16 – Histogram of residuals of mod.Spher and associated QQ plot.

Example Las Rosas (8/8)

```
pred_extended <- predict(modSpher)
head(pred_extended)

##           1            2            3            4            5            6
## 0.26553978 0.20447253 0.14433762 0.09495807 0.06079698 0.05992729

Xutm2 <- Xutm
Xutm2@data$N <- Xutm@data$N + 1
newpred_extended <- predict(modSpher, newdata=Xutm2@data)
head(newpred_extended)

## [1] 0.26756172 0.20649447 0.14635956 0.09698002 0.06281892 0.06194923
```

In practice

- Easier to use a regression model designed for spatially autocorrelated data, and often **more intuitive**.
- If one of these two models does not give a satisfactory result, you can try an extended linear model \Rightarrow choose the form of Λ yourself, using the form of the semi-variogram or criteria like AIC.

Table of Contents

- 1** Origins and Consequences of Spatial Autocorrelation
- 2** Working example : Las Rosas
- 3** Spatial Lag Model
- 4** Spatial Error Model
- 5** Choosing Between Spatial Lag, Error and SAC models
- 6** Extended Linear Models
- 7** Bibliography

Bibliography

- 1** Mixed Effects Models and Extensions in Ecology with R. A. Zuur, E.N. Ieno, N. Walker, A.A. Saveliev and G.M. Smith, Springer, 2009.
- 2** Spatial Data Analysis in Ecology and Agriculture using R. R. E. Plant, CRC Press, 2012.
- 3** Mixed-Effects Models in S and S-PLUS. J.C. Pinheiro and D.M. Bates, Springer, 2000.
- 4** Spatial Processes, Models and Applications. A.D. Cliff and J.K. Ord, Pion Limited, 1981.
- 5** Statistics for Spatial Data, revised edition. N.A.C. Cressie, Wiley Classics Library, 2015.