# Genes mirror geography within Europe

John Novembre,[1,2] Toby Johnson,[4,5,6] Katarzyna Bryc,[7] Zoltán Kutalik,[4,6] Adam R. Boyko,[7] Adam Auton,[7] Amit Indap,[7] Karen S. King,[8] Sven Bergmann,[4,6] Matthew R. Nelson,[8] Matthew Stephens,[2,3] and Carlos D. Bustamante[7]

# Global and local spatial autocorrelation

Dr Stéphane Joost – Oliver Selmoni (Msc)

Laboratory of Geographic Information Systems (LASIG)

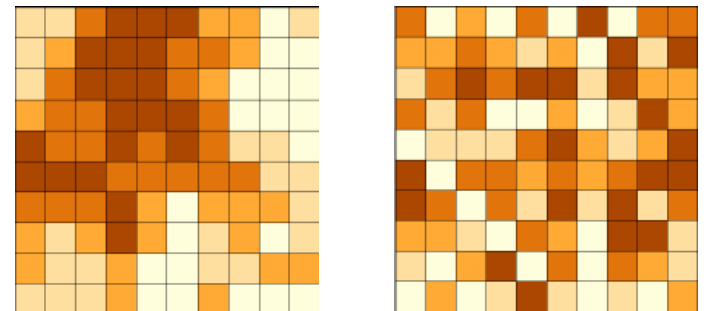Ecole Polytechnique Fédérale de Lausanne (EPFL), Switzerland

# Outline

- Introduction: spatial dependence
- Measuring spatial dependence : Moran's I
- Weighting
- Significance
- Local spatial autocorrelation
- Examples
- Conclusion

# MEASURING SPATIAL DEPENDENCE

- We want to measure how similar are the different values of a given variable for a set of spatially distributed individuals…

  → To quantify the spatial regularity of a given phenomenon

  → To determine the range of spatial dependence
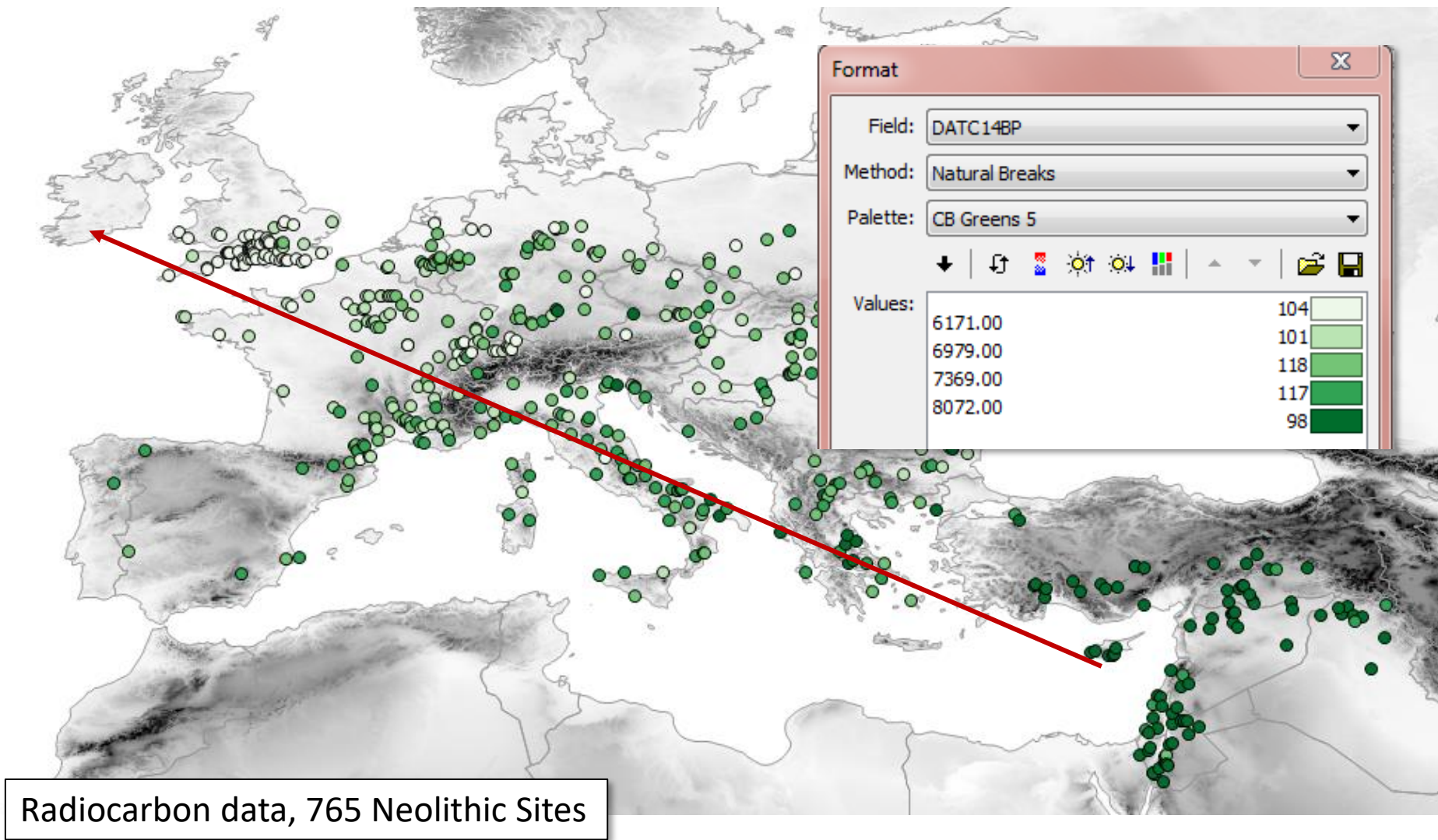
**What is spatial dependence ?**

# Spatial dependence



- A yellow people is more likely to interact with another yellow people
- Similarly, a r&w people is more likely to interact with another r&w people
- The membership determined the spatial distribution of people
- Spatial dependence induced by this membership is perceptible in space through colors

# Tracing the Origin and Spread of Agriculture in Europe
## (Pinhasi et al. 2005, PLOS Biology)



Radiocarbon data, 765 Neolithic Sites

# SPATIAL AUTOCORRELATION: A PARADOX

- Spatial depdendence can be measured by means of indices of **spatial autocorrelation**
- A **paradox**
- First law of geography: "Everything is related to everything else, but near things are more related than distant things", **W.Tobler** (1970)
- This means that natural phenomena (e.g. temperature) as well as socio-demographic ones (e.g. population density) **are not spatially distributed at random**

- **But** to measure the spatial structure of these phenomena, we have to use classical statistical tools **requiring a random spatial distribution of samples** and **independence** between them
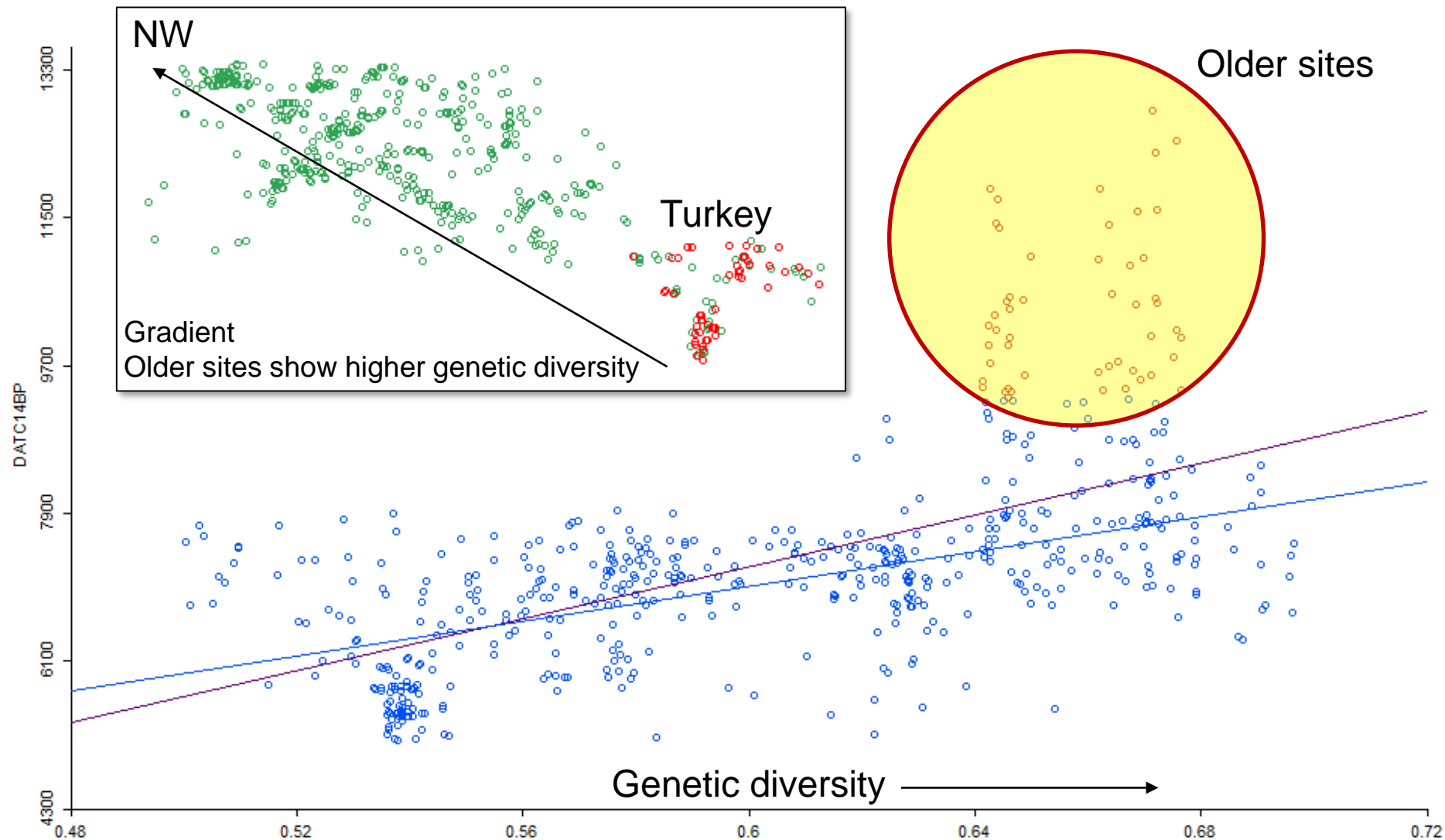
# STANDARD STATISTICS

- Classical statistics are non-spatial

- Based on a neutral geographic space hypothesis

- Geographic space should be the simple **neutral** support for studied phenomena

- Theoretically, the location of a set of observations in space **should not influence their attributes**

# Spatial dependance - autocorrelation

- **Geographical space is not neutral** and consequently many statistical tools are not appropriate
- For instance, ordinary linear regressions (OLR) should be implemented **only if observations are selected at random**
- When observations show spatial dependence, estimated values for the whole data set **are distorted/biased**
- Indeed, sub-regions including a concentration of individuals with high values will have an important impact on the model and **lead to a global overestimation** of the variable under study over the whole area
- In other words, a strong correlation between two variables **at a single location** of the territory will influence the measure of the relationship **over the whole territory**

NW

Turkey

Gradient
Older sites show higher genetic diversity

Older sites

Genetic diversity ⟶

DATC14BP

microsatOb

| #obs | R^2 | const a | std-err a | t-stat a | p-value a | slope b | std-err b | t-stat b | p-value b |
|------|-------|-----------|-----------|----------|-----------|-----------|-----------|----------|-----------|
| 538 | 0.368 | -2.24e+003 | 544 | -4.12 | 4.38e-005 | 1.58e+004 | 895 | 17.7 | 0 |
| 488 | 0.355 | 653 | 390 | 1.67 | 0.0948 | 1.06e+004 | 649 | 16.4 | 0 |

We have to use standard statistical tools with caution with spatial data (remember this paradox)

# INDICES TO MEASURE SPATIAL AUTOCORRELATION
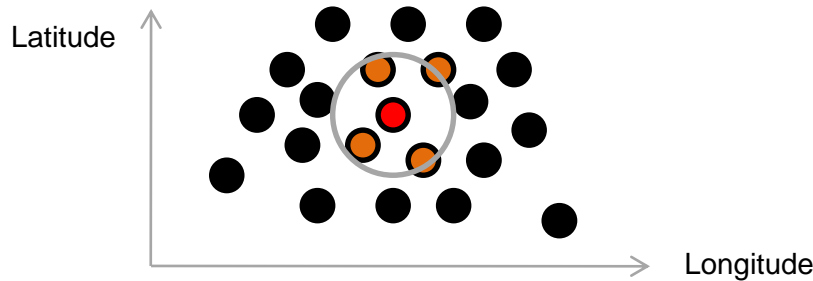
Developed with standard statistics

- Geary's C
- Ripley's K
- Join Count Analysis
- **Moran's I**

# MEASURING SPATIAL AUTOCORRELATION

- Spatial autocorrelation indices express **the degree of spatial structuring** of a given variable

- Spatial autocorrelation indices permit to **quantify the spatial regularity** of a geographically distributed phenomenon

- Spatial autocorrelation is **positive** when values measured at neighbouring points **resemble each other**

- It is **negative** in case of dissimilarity

# **Neighborhood relationship** and **spatial weighting**

- We want to know how similar is the value *v* of object *o* **in comparison** with the value *v* of other objects located in its neighbourhood ? (to measure spatial dependence)

Latitude

Longitude

- We need to define this **neighbourhood**: it may be a distance (100m around each object) or a given number of neighbours (4 nearest neighbours for instance)
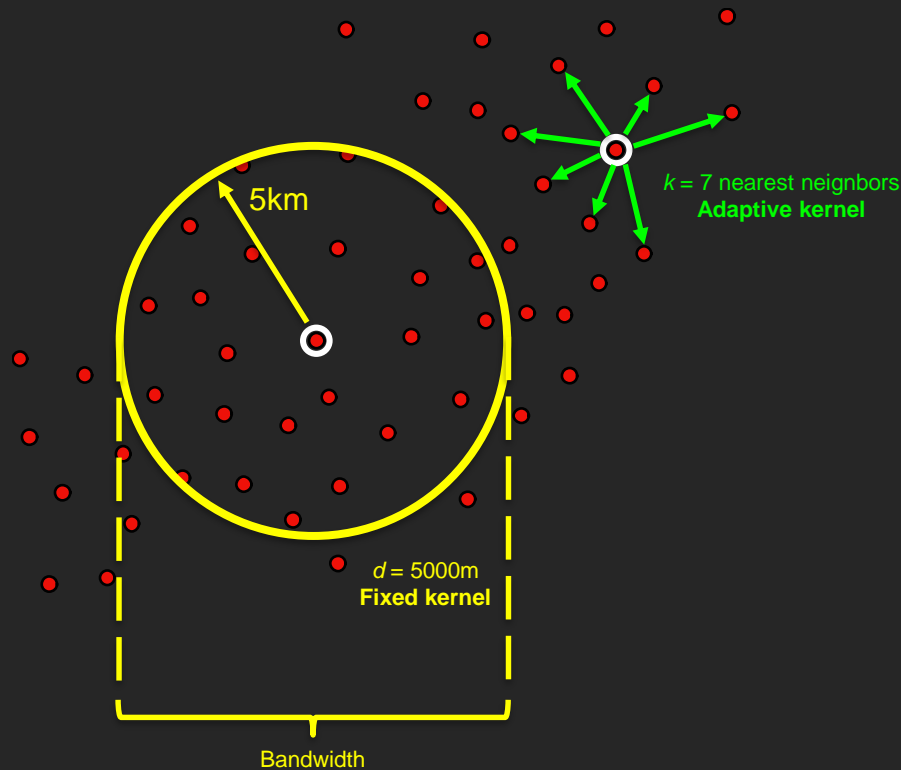
- This criterion (100m or 4 nearest neighbours) defines the spatial weight

# Neighborhood relationships

- Spatial autocorrelation is characterized by a correlation between measures of a given phenomenon located close to each other



5km

Etc.

To quantify the spatial dependence and produce a measure of global spatial autocorrelation, it is necessary to take into account the neigborhood of each of the considered geographic objects

# Several ways to define the spatial neighborhood

# How to define the neighborhood of polygons

# Spatial weighting

**Spatial weighting scheme**

Fixed kernel $\longrightarrow$ Bandwidth $d$

Adaptive kernel $\longrightarrow$ $k$ nearest neighbors

Contig order $n$

Spatial weighting file

```
0 54 comvd_prec ide
1 2       1895.70644
1 13      2031.2413
1 4       2062.65071
1 35      2365.3336
1 3       2474.90419
1 29      2533.8499
1 19      3041.1263
...
```

```
comvd_prec_qu1.gal
1    0 54 comvd_prec ide
2    1 6
3    35 29 13 4 3 2
4    2 5
5    9 8 4 3 1
6    3 5
7    13 9 21 1 2
8    4 6
9    8 7 5 1 2 29
10   5 2
11   7 4
12   6 2
13   10 9
14   7 3
15   8 4 5
16   8 5
17   10 9 2 4 7
18   9 5
19   10 3 2 6 8
20   10 3
21   6 8 9
22   11 8
23   35 27 26 21 15 13 14 18
```

$n = 1$

```
comvd_prec_7k.gwt
1    0 54 comvd_prec ide
2    1 2       1895.70644
3    1 13      2031.24132
4    1 4       2062.65071
5    1 35      2365.33363
6    1 3       2474.90419
7    1 29      2533.84993
8    1 19      3041.12639
9    2 3       1813.98042
10   2 1       1895.70644
11   2 9       1992.83081
12   2 8       2039.59784
13   2 4       2479.38678
14   2 7       2693.1369
15   2 13      3076.66614
```

$k = 7$

29
35
4
1
13
2
3

# MORAN'S I

- Moran's autocorrelation coefficient I is an extension of Pearson product-moment correlation coefficient

$$I = \frac{N \sum_i \sum_j W_{i,j}(X_i - \bar{X})(X_j - \bar{X})}{(\sum_i \sum_j W_{i,j}) \sum_i (X_i - \bar{X})^2}$$

Where

- **N** is the number of observation units
- **W**$_{ij}$ is a spatial weight applied to define the comparison between locations **i** and **j**
- **X**$_i$ is the value of the variable at a location i
- **X**$_j$ is the value of the variable at a location j
- **X̄** is the mean of the variable
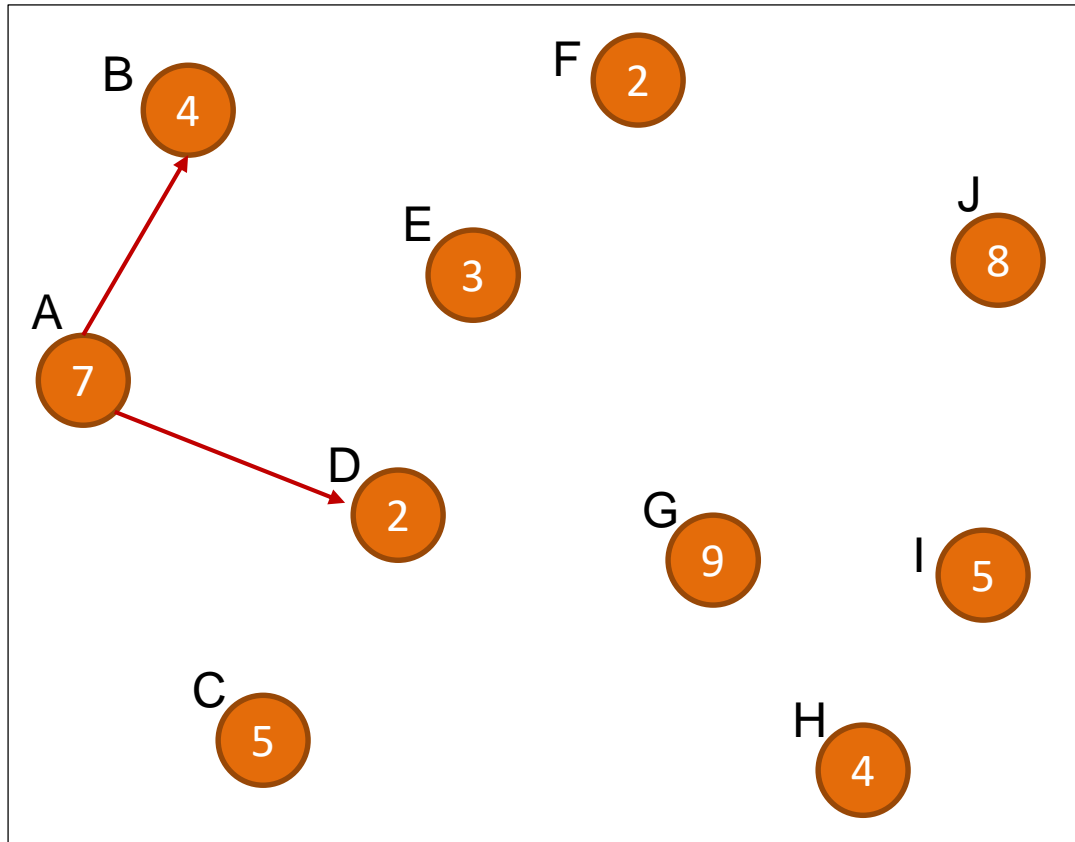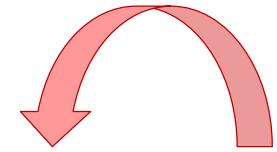
# Moran's I as a regression coefficient

- Anselin (1996): Moran's I can be interpreted as a regression coefficient

- $\rightarrow$ regression of *WZ* on *Z*

- The *weighted value of Z* (mean of *Z* according to the weighting criteria) **on** *Z* (the observed value at the point of interest)

- This interpretation provides a way to visualize the linear association between *Z* and *WZ* in the form of a bivariate scatterplot

- Anselin (1996) referred to this plot as the Moran scatterplot

- He pointed out that **the least squares slope** in a regression through the origin **is equal** to **Moran's I**

# PROCESSING MORAN'S I

Weighting criteria: 2 nearest neighbours
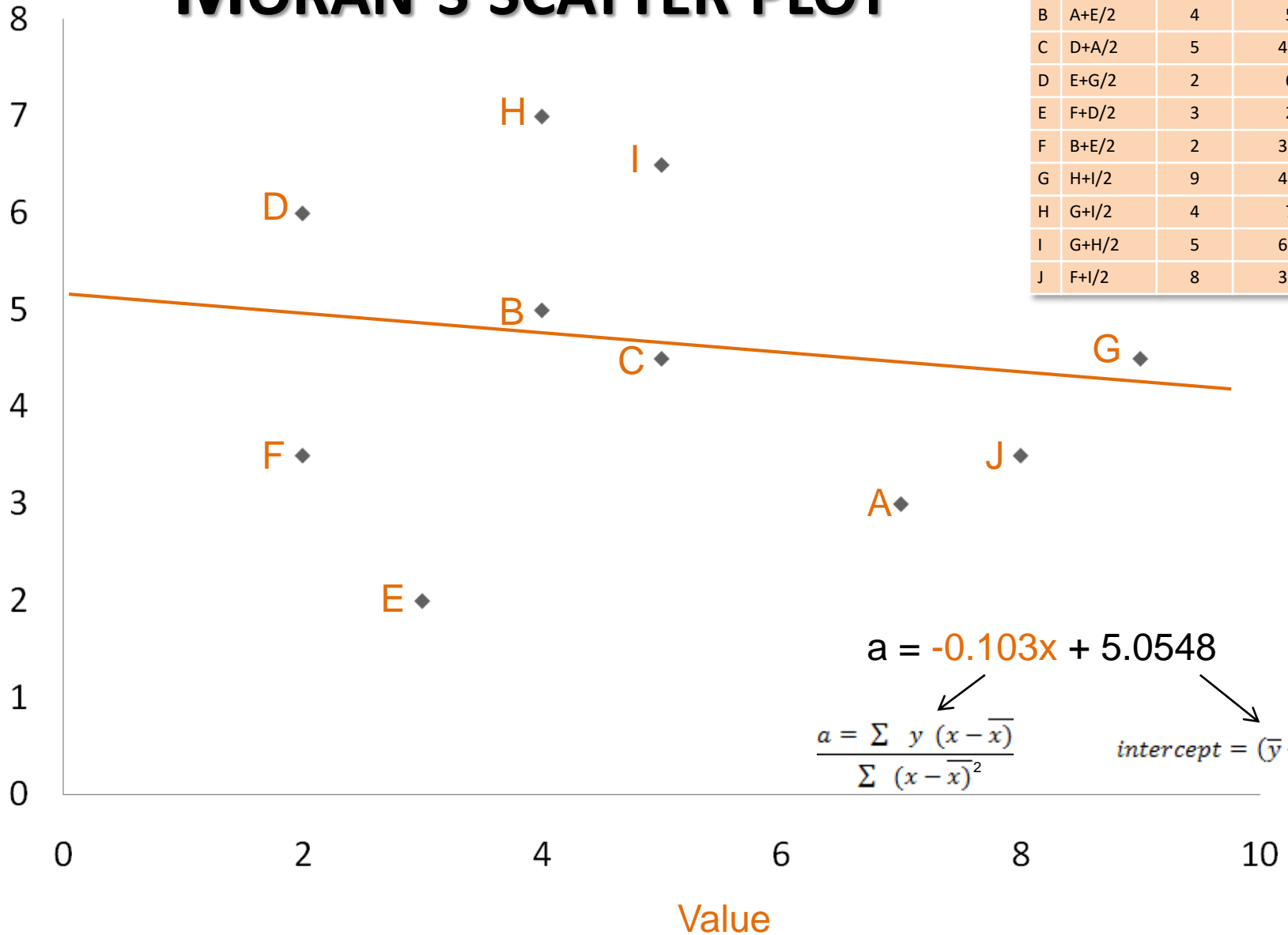The value of the variable is displayed in the centroids

Regression of Y on X

|  | Weighting | Value (X) | Mean of nearest neighbours values (Y) |
|---|---|---|---|
| A | B+D/2 | 7 | 3 |
| B | A+E/2 | 4 | 5 |
| C | D+A/2 | 5 | 4.5 |
| D | E+G/2 | 2 | 6 |
| E | F+D/2 | 3 | 2 |
| F | B+E/2 | 2 | 3.5 |
| G | H+I/2 | 9 | 4.5 |
| H | G+I/2 | 4 | 7 |
| I | G+H/2 | 5 | 6.5 |
| J | F+I/2 | 8 | 3.5 |

# MORAN'S SCATTER PLOT

| | Weighting | Value | Mean of NN |
|---|---|---|---|
| A | B+D/2 | 7 | 3 |
| B | A+E/2 | 4 | 5 |
| C | D+A/2 | 5 | 4.5 |
| D | E+G/2 | 2 | 6 |
| E | F+D/2 | 3 | 2 |
| F | B+E/2 | 2 | 3.5 |
| G | H+I/2 | 9 | 4.5 |
| H | G+I/2 | 4 | 7 |
| I | G+H/2 | 5 | 6.5 |
| J | F+I/2 | 8 | 3.5 |

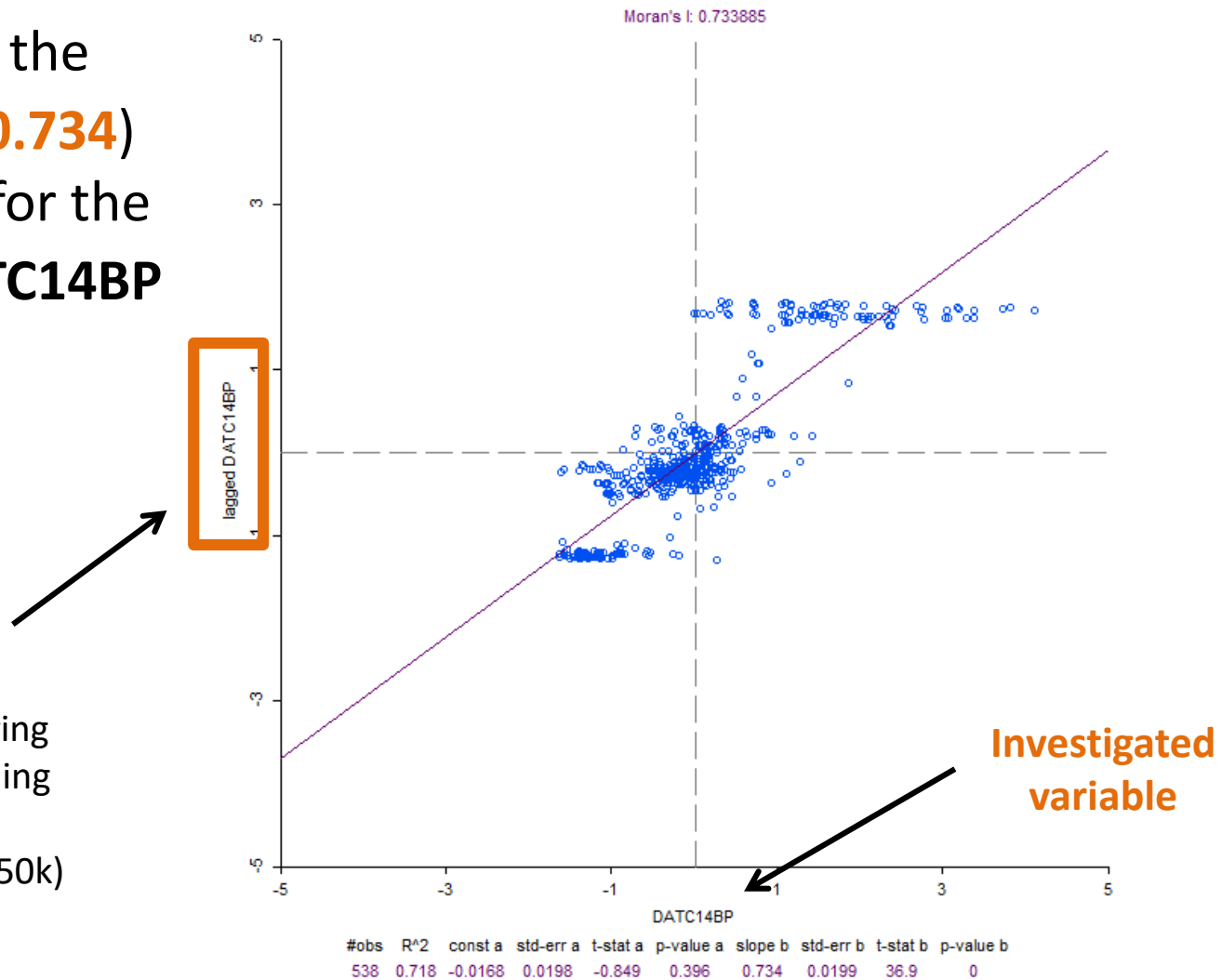Weighted value (mean of Nearest Neighbours)

Value

$a = -0.103x + 5.0548$

$$a = \frac{\sum y\,(x - \overline{x})}{\sum (x - \overline{x})^2}$$

$$intercept = (\overline{y} - a)\,\overline{x}$$

# I = SLOPE OF THE REGRESSION

The slope of the regression (**0.734**) is Moran's I for the Variable **DATC14BP**

**Weighted variable**: mean of neighbouring spatial units according to the selected weighting criteria (50k)



**Investigated variable**

| #obs | R^2 | const a | std-err a | t-stat a | p-value a | slope b | std-err b | t-stat b | p-value b |
|------|-----|---------|-----------|----------|-----------|---------|-----------|----------|-----------|
| 538 | 0.718 | -0.0168 | 0.0198 | -0.849 | 0.396 | 0.734 | 0.0199 | 36.9 | 0 |

# Moran's I range of values

- Moran's I statistics ranges **from -1 to 1**
- A value close to 1 shows a strong positive spatial autocorrelation
- A value close to -1 shows a strong negative spatial autocorrelation (opposition between individuals)
- **0** = no spatial autocorrelation **= independence between individuals** = neutral geographic space
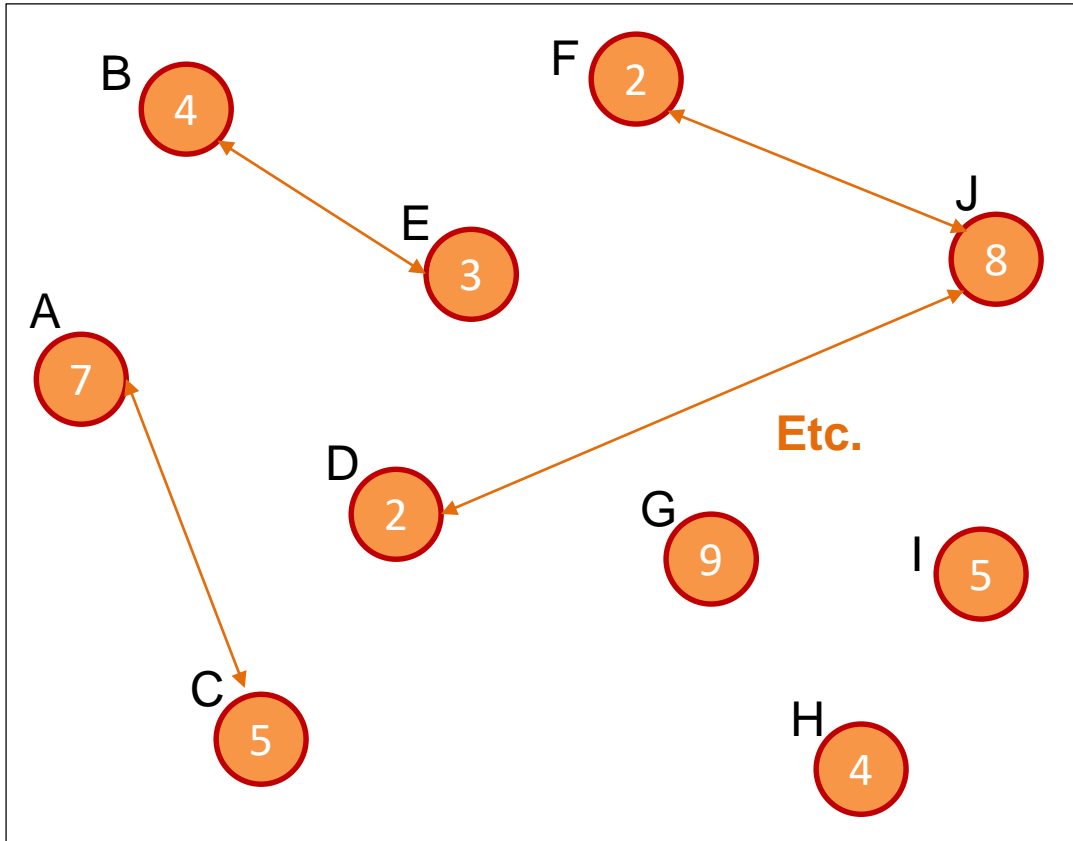
# MORAN'S I SIGNIFICANCE

- We need a statistical test to determine whether data attached to individuals are distributed at random over the territory

- We use random **permutations** between individuals to perform this test

- How does the observed situation behave **in comparison with all other possible configurations** ?

# Random permutations

- For each run, the **attributes** of all individuals in the data set **are randomly moved** between **all** possible **locations**
- Many runs are performed by means of Monte-Carlo method (e.g. 500, 1000, or more permutations)
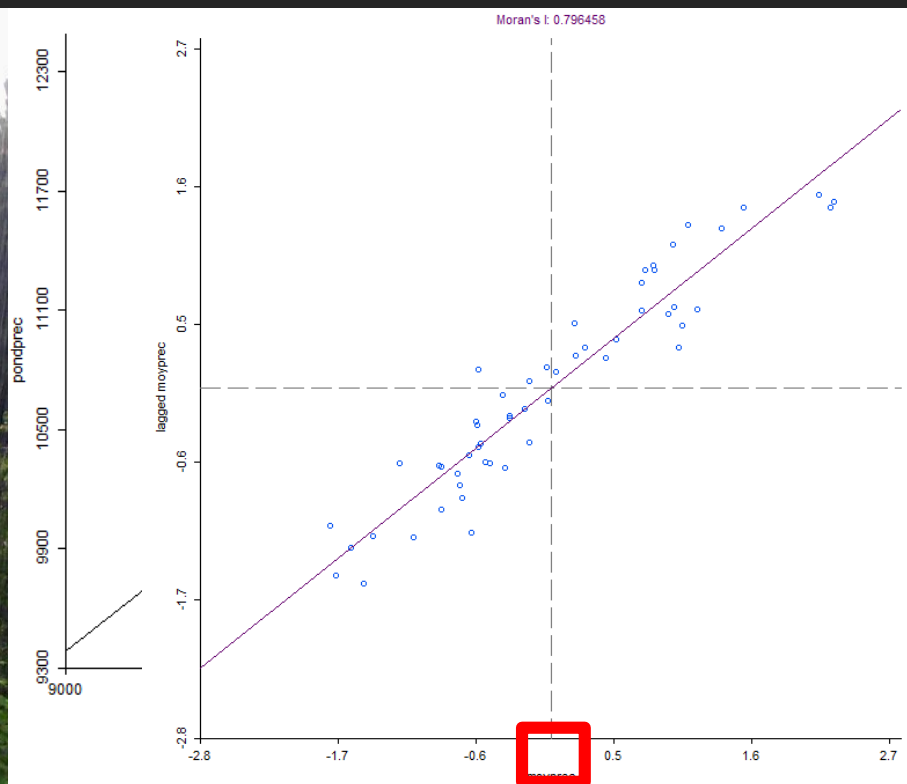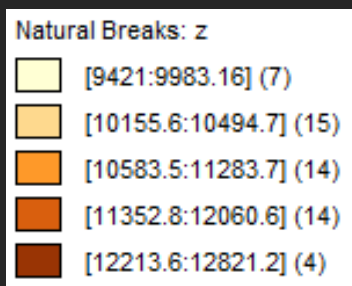- More runs = more significance

# RANDOM PERMUTATIONS



- Moran's I is **calculated for each run**

- The total number of possible permutations is **n!** (here 10! = 3'628'800)

# Moran's I as a coefficient of regression

# Visualization of the spatial structure



Natural Breaks: z
- [9421:9983.16] (7)
- [10155.6:10494.7] (15)
- [10583.5:11283.7] (14)
- [11352.8:12060.6] (14)
- [12213.6:12821.2] (4)

I = 0.79

# Moran's I significance and random permutations

Observed situation : $I = I_0$

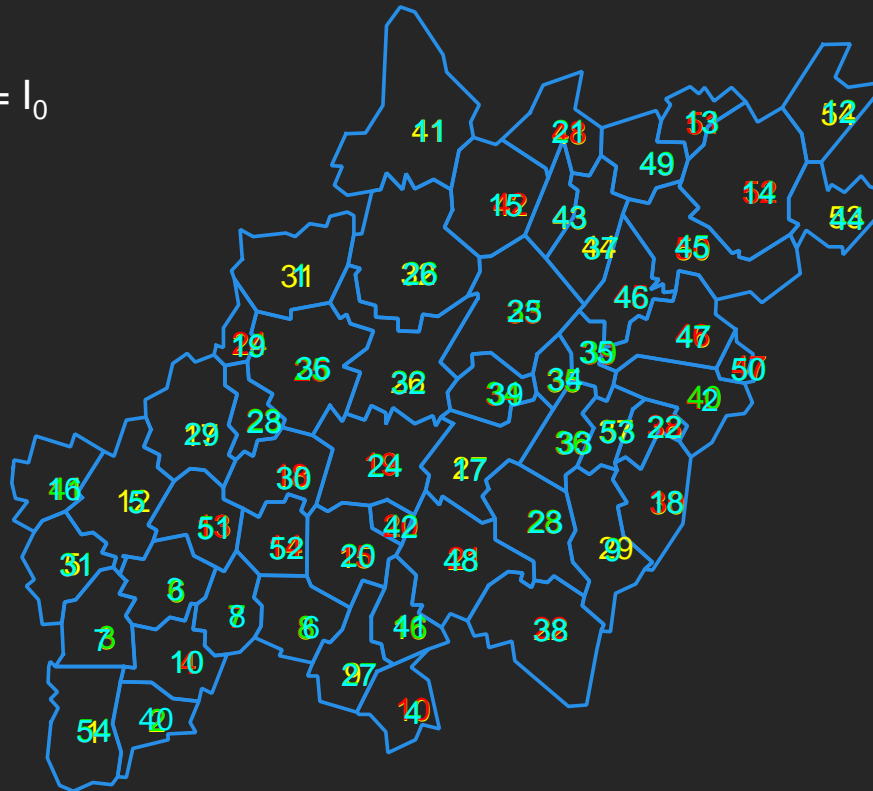Draw 1: $I = I_1$

Draw 2: $I = I_2$
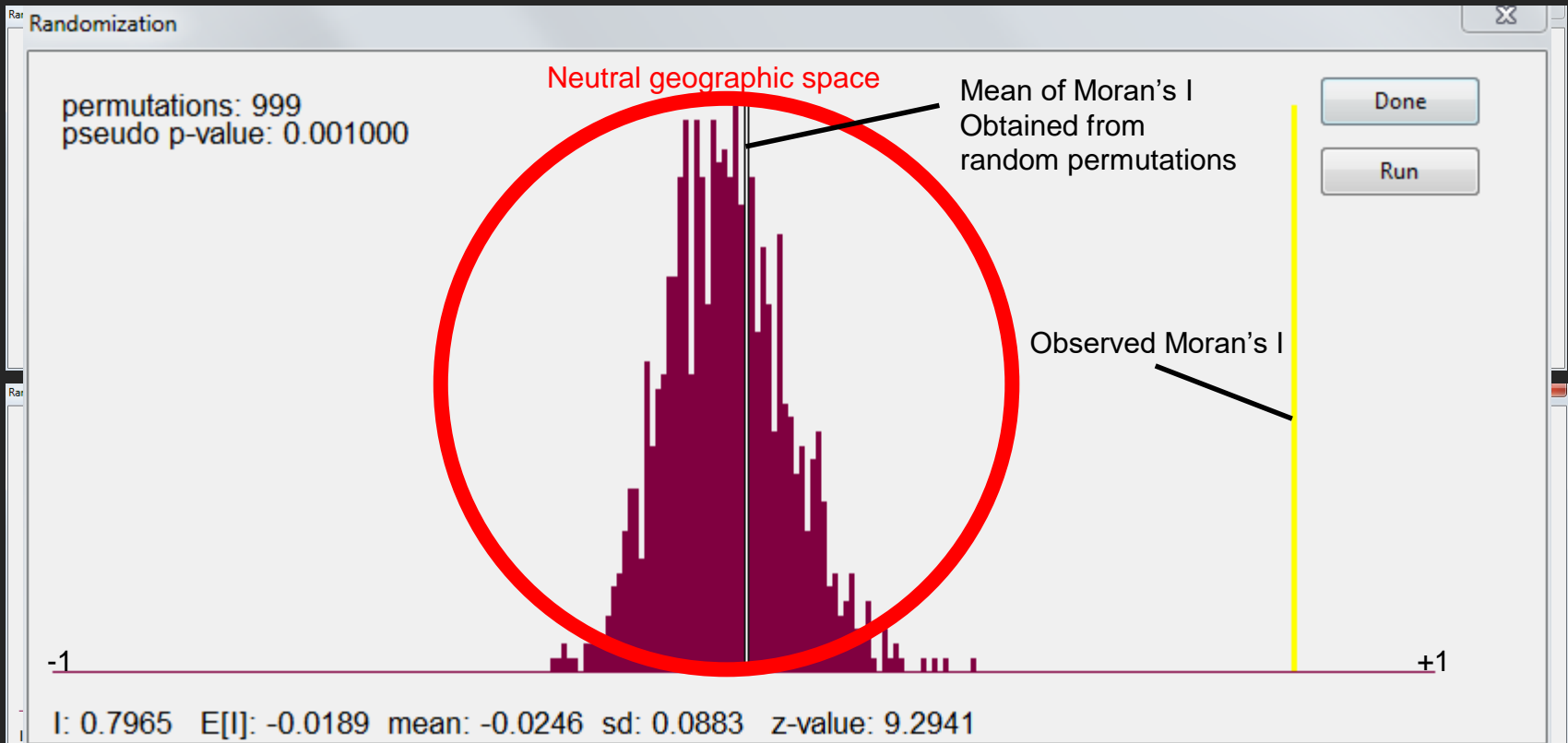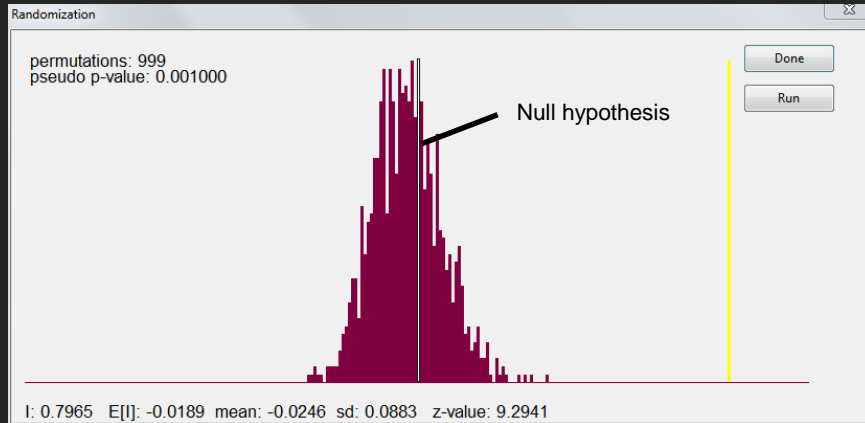
Draw 3: $I = I_3$

Draw 4: $I = I_4$

...

...

54! possible configurations

# Histogram of random  situations and p-values

# How to calculate the significance and the p-value



permutations: 999
pseudo p-value: 0.001000

Null hypothesis

Done

Run

I: 0.7965  E[I]: -0.0189  mean: -0.0246  sd: 0.0883  z-value: 9.2941

$$\text{p-value} = \frac{0+1}{999+1} = 0.001$$

$$\text{p-value} = \frac{Nb\ I_{al} \geq I_{obs} + 1}{Nb\ permutations + 1}$$

$$ou\quad \frac{Nb\ I_{al} \leq I_{obs} + 1}{Nb\ permutations + 1}$$

A Moran's I of 0.79 translates a spatial structure significantly different from a random spatial distribution

# LISA : Local Indicators of Spatial Association

- It is also possible to calculate **local Moran's indices**

- Moran's I is decomposed into several local coefficients, whose sum over the whole studied area is proportional to the global Moran's I

- The significance is assessed locally (same procedure like Global I except that the value of the point of interest remains fixed, and neighbours are among N-1 values)

# Calculation of Local I

| 45 | 44 | 44 |
|----|----|----|
| 43 | 42 | 39 |
| 38 | 32 | 34 |

Rook weighting scheme

| $y_i$ | $z_i$ | $w_{ij}$ | $w_{ij}z_j$ |
|-------|-------|----------|-------------|
| 45 | 4.889 | 0 | 0 |
| 43 | 2.889 | 0.25 | 0.722 |
| 38 | -2.111 | 0 | 0 |
| 44 | 3.889 | 0.25 | 0.972 |
| 42 | 1.889 | 0 | 0 |
| 32 | -8.111 | 0.25 | -2.028 |
| 44 | 3.889 | 0 | 0 |
| 39 | -1.111 | 0.25 | -0.278 |
| 34 | -6.111 | 0 | 0 |
| | | 1 | -0.611 |

$$I_i = \left[\frac{Z_i}{S^2}\right] \sum_{j=1}^{n} w_{ij}z_j , j \neq i$$

y = value
Z = dev. from the mean (40.1)
w = weight

Mean = 40.1
$S^2$ = Variance = 21.8 $\quad \frac{\sum (x - \bar{x})^2}{(n-1)}$

- Here weights are standardized (they sum to 1). We have 4 values (because of the Rook criterion).
- **For this location**, $z_i$ = 42 - 40.111 = 1.889
- The sum of the weights multiplied by the deviations from the mean = -0.611
- $I_i$ = 1.889/21.861 x -0.611 = -0.053
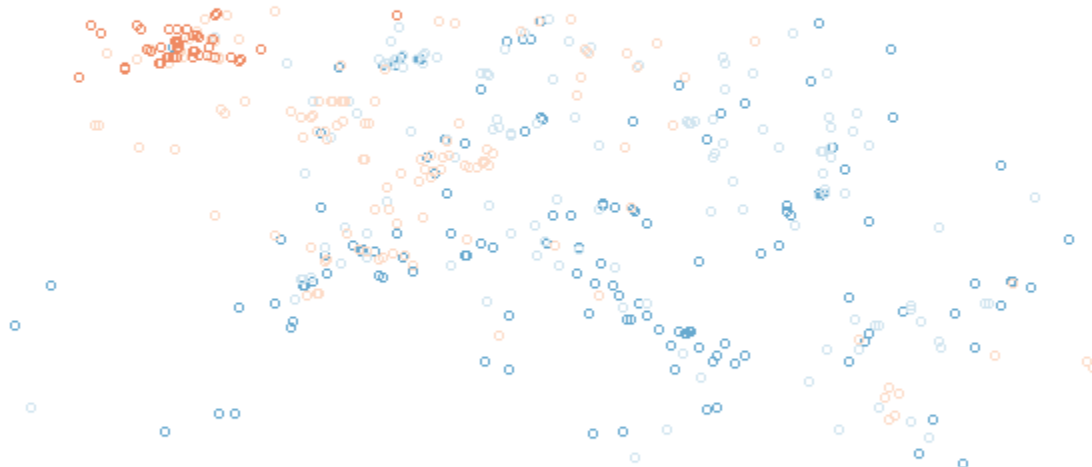- Local I values sum up to global Moran's I

# LISA

Thus it is possible:

– to map indices and corresponding p-values to show how local spatial autocorrelation varies over the territory

– to highlight local regimes of spatial autocorrelation
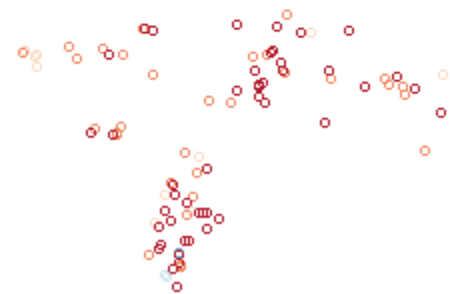
# Local Index of Spatial Association (DATC14BP)



Hinge=1.5: LISA_I_DAT

- Lower outlier (0)
- < 25% (134)
- 25% - 50% (135)
- 50% - 75% (136)
- > 75% (86)
- Upper outlier (47)

**Table**

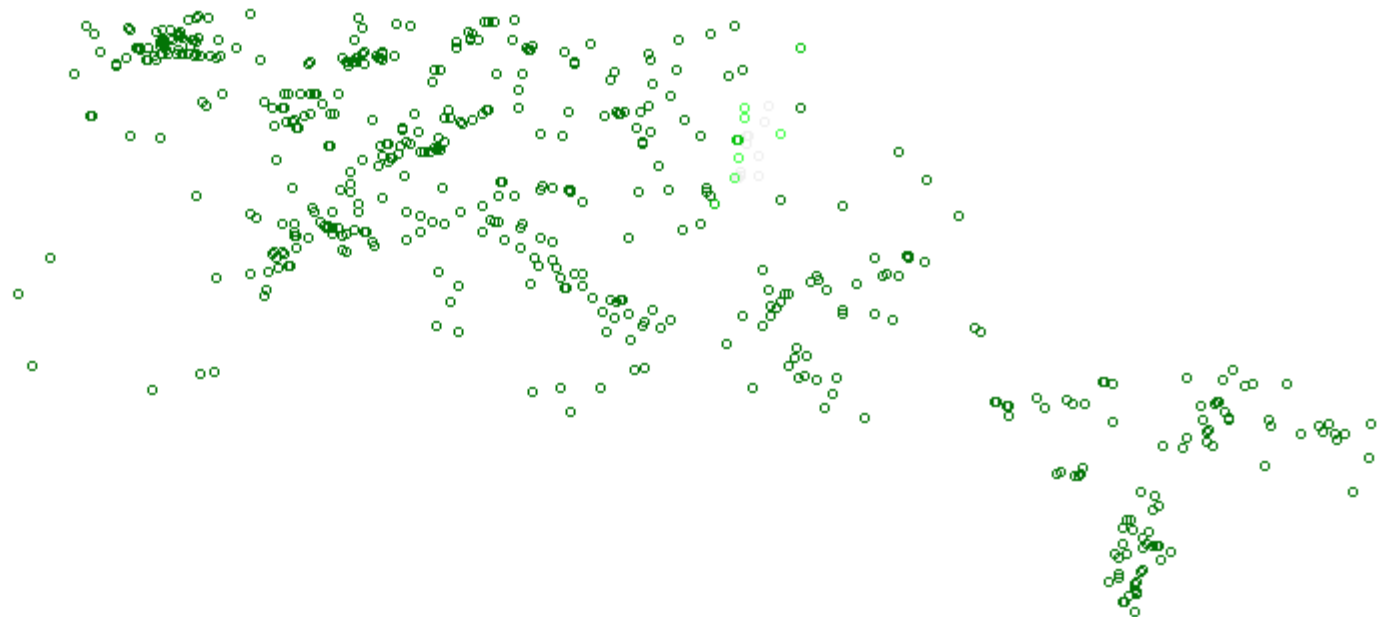| | LISA_I_DAT < |
|---|---|
| 6 | 4.5490719 |
| 8 | 3.6584164 |
| 70 | 3.3732070 |
| 4 | 3.2160546 |
| 67 | 2.8434214 |
| 27 | 2.6163152 |
| 2 | 2.5946365 |
| 50 | 2.5303988 |
| 16 | 2.4112840 |
| 31 | 2.3624535 |
| 49 | 2.2583355 |
| 15 | 2.1428783 |
| 52 | 2.0865199 |
| 66 | 2.0524817 |
| 44 | 2.0123931 |

Min = - 0.36

# Significance map



LISA Significance Map: sitesNeolithique25km,

- ☐ Not Significant (11)
- ☐ p = 0.05 (7)
- ☐ p = 0.01 (4)
- ☐ p = 0.001 (516)
- ☐ p = 0.0001 (0)

9999 permutations

# LISA CLUSTERS



Moran's I: 0.323921

LOW-HIGH                    HIGH-HIGH

lagged DATC14BP

LOW-LOW                     HIGH-LOW

DATC14BP

| #obs | R^2 | const a | std-err a | t-stat a | p-value a | slope b | std-err b | t-stat b | p-value b |
|------|-----|---------|-----------|----------|-----------|---------|-----------|----------|-----------|
| 538 | 0.612 | -0.0658 | 0.0111 | -5.91 | 5.99e-009 | 0.324 | 0.0111 | 29.1 | 0 |

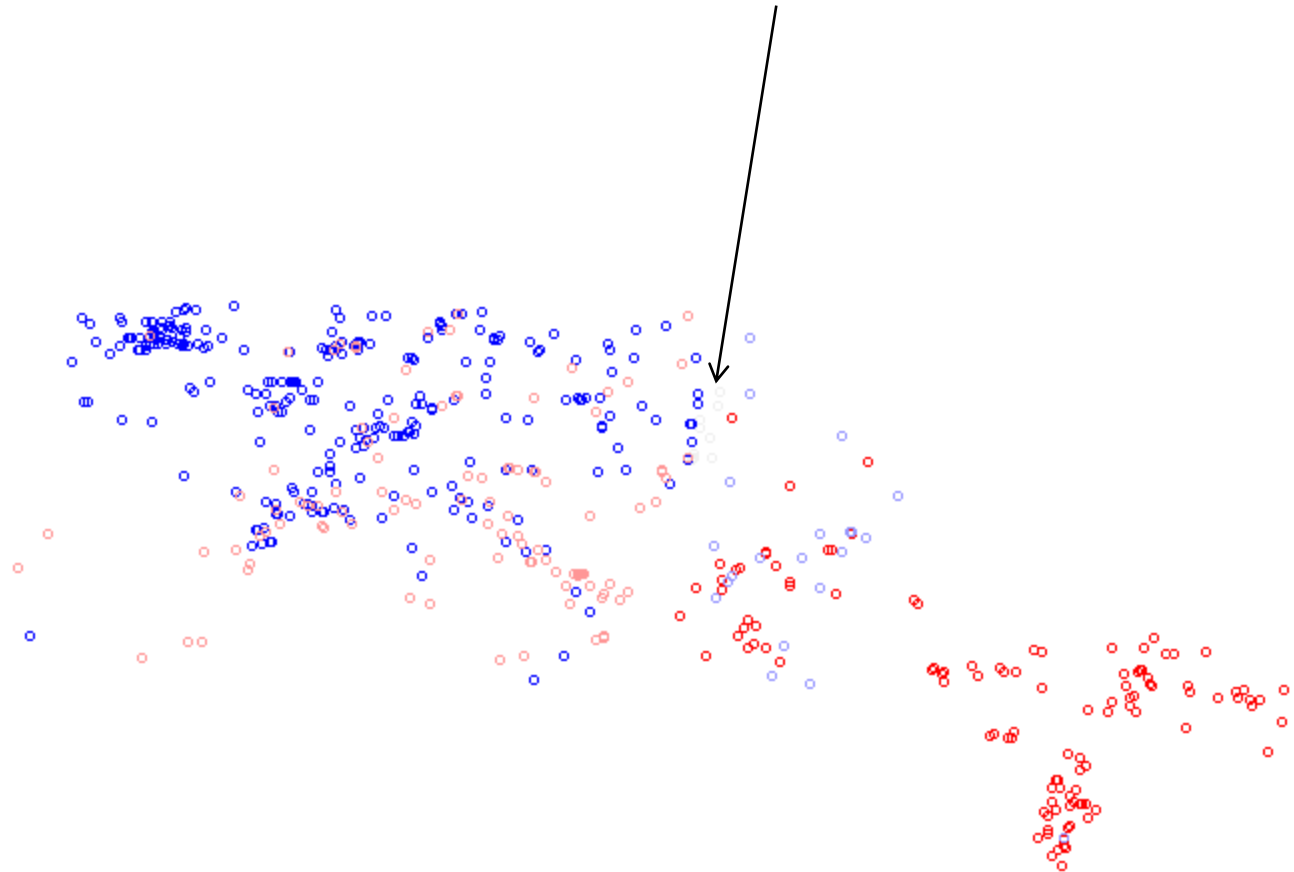# Local spatial autocorrelation on DATC14BP



LISA Cluster Map: sitesNeolithique25km, I_DATC14BP

- ☐ Not Significant (11)
- ■ High-High (120)
- ■ Low-Low (277)
- ■ Low-High (21)
- ■ High-Low (109)

# IN SUMMARY

- To assess spatial dependence:
  - Calculate global or local indicators of spatial autocorrelation
  - Indicators based on the resemblance between points of interest and their neighborhood
  - Significance: establish a comparison with a **neutral geographic space** (null hypothesis = the spatial distribution of attributes is random)
- We calculate spatial dependence with theories relying on contradictory arguments (space is not neutral vs space must be neutral)
- → **Geographically Weighted Regression** (GWR)
- Includes spatial weighting in the processing of the regression

# REFERENCES

- Anselin, L. (1995). Local indicators of spatial association — LISA. Geographical Analysis, 27:93–115.
- Anselin, L. (1996). The Moran scatterplot as an ESDA tool to assess local instability in spatial association. In Fischer, M., Scholten, H., and Unwin, D., editors, Spatial Analytical Perspectives on GIS in Environmental and Socio-Economic Sciences, pages 111–125. Taylor and Francis, London.
- Legendre, P. (1993) Spatial Autocorrelation: Trouble or New Paradigm? Ecology, Vol. 74, No. 6, pages 1659-1673