

Identifying urban profiles with regard to overheating factors

Application to French municipalities

Pauline Besnard, Baptiste Alglave, Jérémie Bernard, Matthieu Goussef, François Leconte, Erwan Bocher

10 ÈMES RENCONTRES DE STATISTIQUE

**SCIENCE
DES DONNÉES
HISTOIRE
& TERRITOIRES**

27 & 28

NOVEMBRE 2025

Amphithéâtre Yves Coppens
Faculté Sciences & Sciences
de l'Ingénieur
Université Bretagne Sud
Campus de Tohannic - VANNES

SCIENTIFIC CONTEXT

- Urban territories are increasingly exposed to **overheating** processes [Masson, 2000].
 - ➡ Threats on populations [Huang et al., 2023]
- Classifying territories with respect to their **sensitivity to overheating** is a critical issue [Wang et al., 2020].
- Still, very few **classifications** with regard to overheating factors.

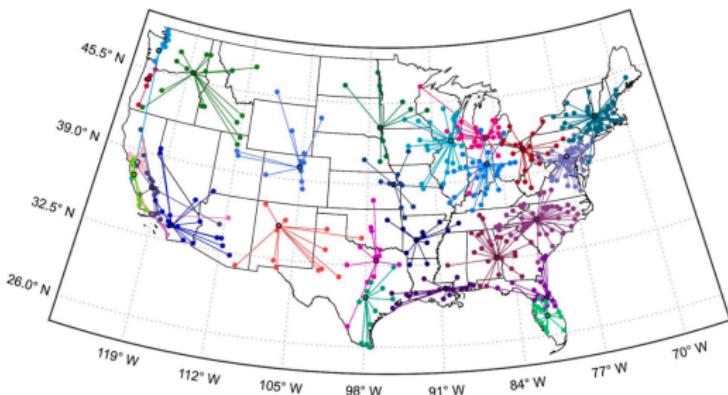
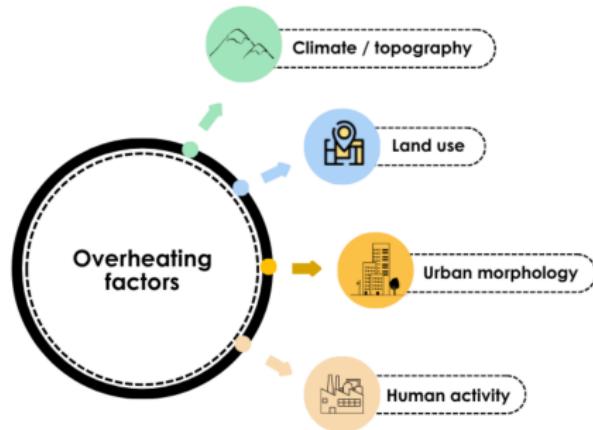
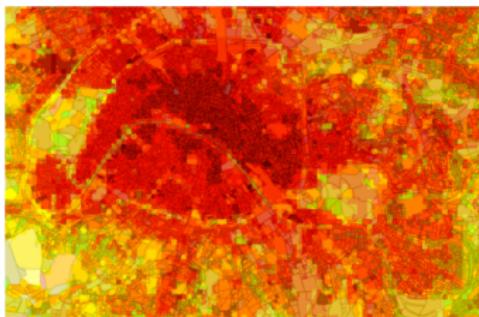
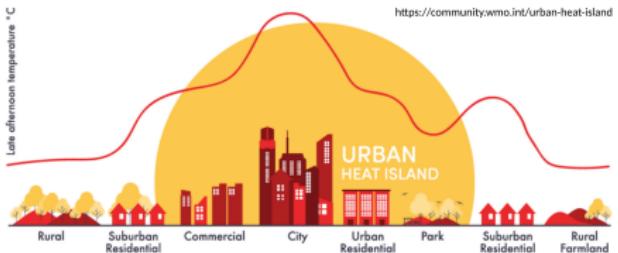
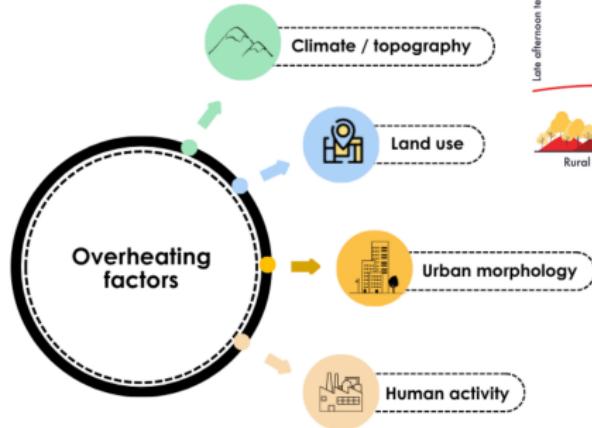


Figure 1: Urban clustering based on geographical distance and 8-day composite daytime LST during a heat wave (July 12–19, 2006).

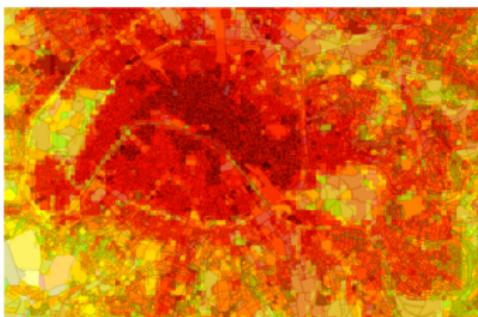
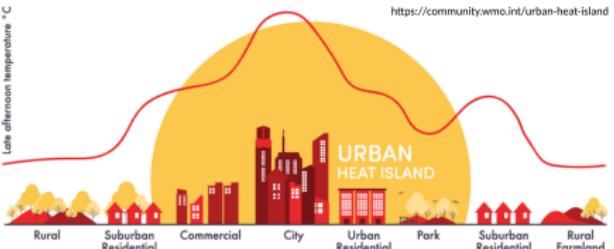
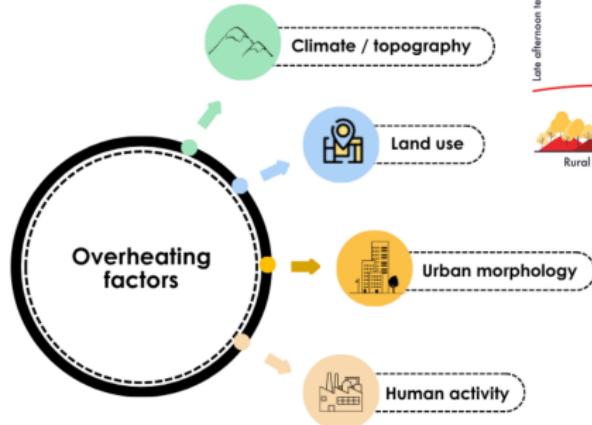
URBAN OVERHEATING FACTORS



URBAN OVERHEATING FACTORS



URBAN OVERHEATING FACTORS



➡ Data to analyze urban areas and inform spatial planning?

POTENTIAL OF GIS TO ANALYZE URBAN AREAS

GIS (Geographic Information System)

- store and access massive and comprehensive amount of geographic data to map territories

BD Topo

Database of the elements of the French territory and its infrastructures (administrative boundaries, buildings, hydrography, land use, transport).



INSTITUT NATIONAL
DE L'INFORMATION
GÉOGRAPHIQUE
ET FORESTIÈRE

OpenStreetMap

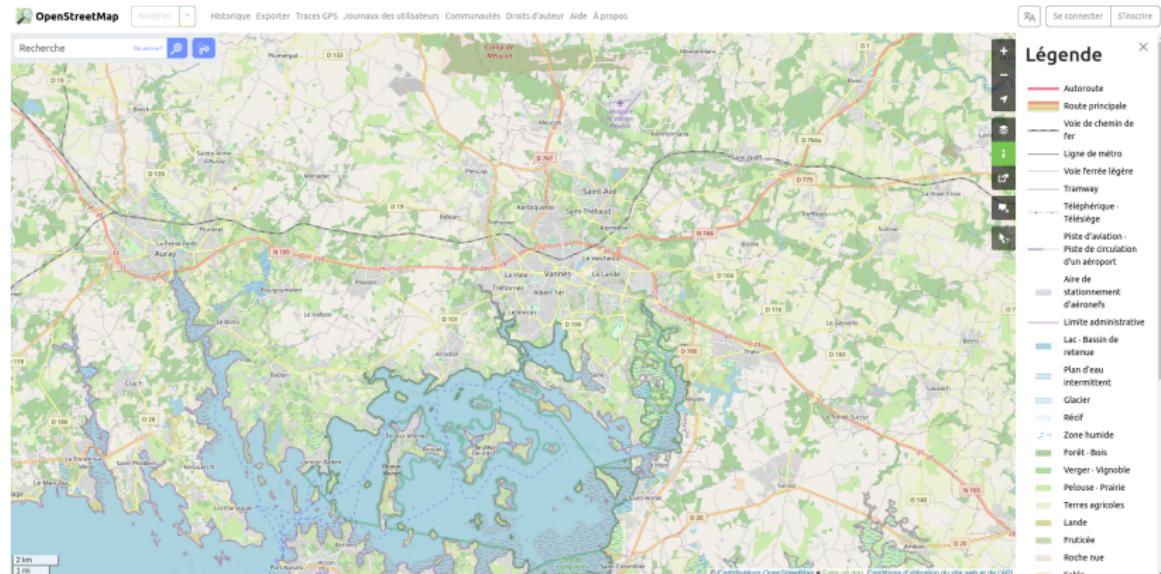
Free, participative, global database
[Vargas-Munoz et al., 2020]



POTENTIAL OF GIS TO ANALYZE URBAN AREAS

GIS (Geographic Information System)

- store and access massive and comprehensive amount of geographic data to map territories



GIS APPLIED TO URBAN CLIMATE

In recent years, GIS tools have been developed to describe and document territories with regard to urban microclimate.

Geoclimate
[Bocher et al., 2021]



Wudapt
[Ching et al., 2018]



Zones Climatiques Locales
[Cerema, 2024]



GIS APPLIED TO URBAN CLIMATE

In recent years, GIS tools have been developed to describe and document territories with regard to urban microclimate.

Geoclimate
[Bocher et al., 2021]



Wudapt
[Ching et al., 2018]



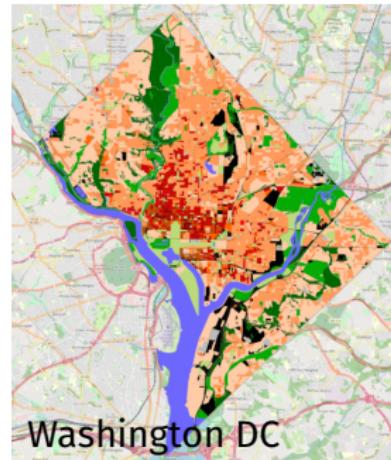
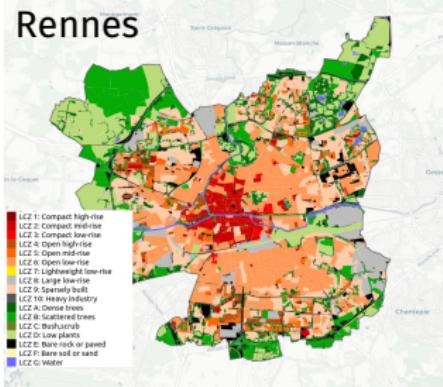
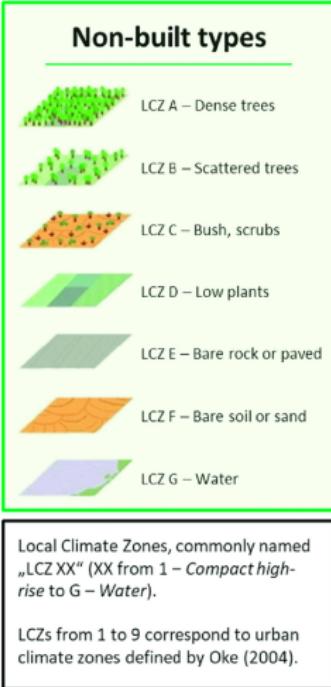
Zones Climatiques Locales
[Cerema, 2024]



- ➡ Rely on the concept of Local Climate Zone (LCZ)

LOCAL CLIMATE ZONES (LCZ)

LCZ types



Aim of this work: propose a methodological framework for classifying urban areas on the basis of urban overheating indicators.

Data of applications: French municipalities

Method:

- Identify overheating factors → define indicators of overheating based on geographic data
- Factorial Analysis of Mixed Data (FAMD) for analysing the indicators combined with k-means for the clustering
- Compare the clusters compared with Urban Heat Island (UHI) model outputs.

SCIENTIFIC ISSUE

Aim of this work: propose a methodological framework for classifying urban areas on the basis of urban overheating indicators.

Data of applications: French municipalities

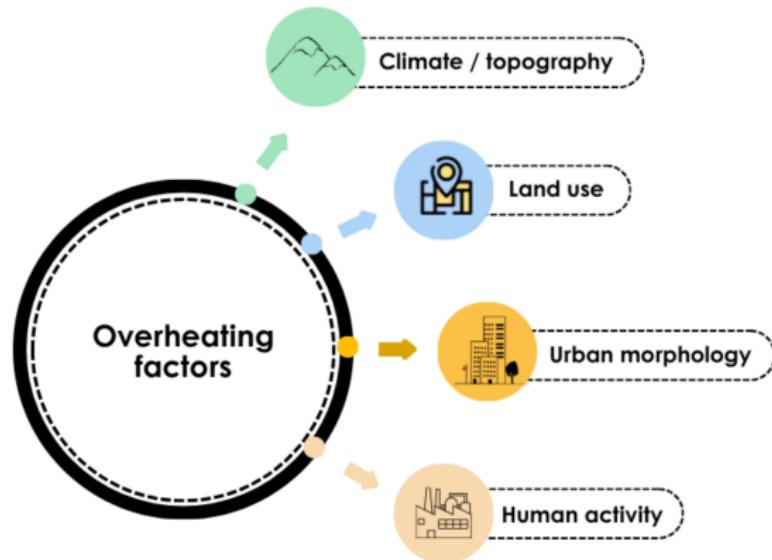
Method:

- Identify overheating factors → define indicators of overheating based on geographic data
- Factorial Analysis of Mixed Data (FAMD) for analysing the indicators combined with k-means for the clustering
- Compare the clusters compared with Urban Heat Island (UHI) model outputs.

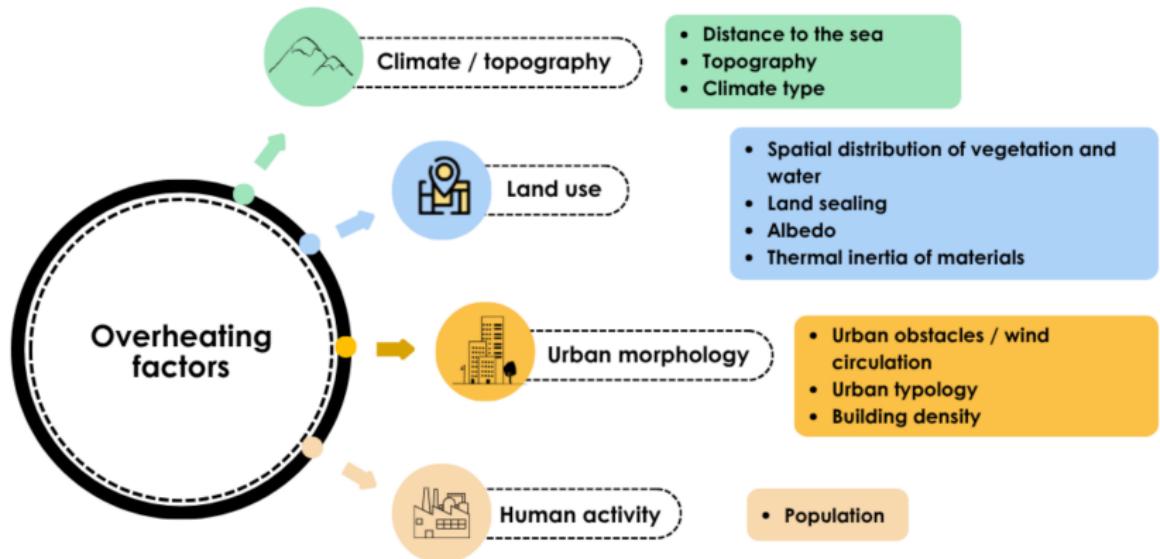


+ additional datasets

URBAN OVERHEATING FACTORS AND INDICATOR DEFINITION



URBAN OVERHEATING FACTORS AND INDICATOR DEFINITION



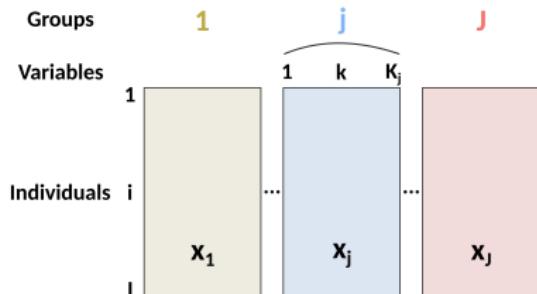
URBAN OVERHEATING FACTORS AND INDICATOR DEFINITION

► 13 indicators describing overheating factors

Table 1: Variables used in the factorial analysis of mixed data

| Class of factor | Variables | Data source (resolution) | Variable type |
|-------------------------------|--|--------------------------|---------------------|
| <i>Climate and topography</i> | Distance of the municipality centroid to the sea | NaturalEarth (1:10e6) | Positive continuous |
| | Climate type (based on Köppen classification) | kgc package (0.05°) | Categorical |
| | Topographic Rugosity Index (TRI) | SRTM (30m) | Positive continuous |
| <i>Land use</i> | Proportion of impervious surface | OSM (Geoclimate) | Proportion |
| | Mean distance between closest vegetation areas | OSM (Geoclimate) | Positive continuous |
| | Proportion of vegetation | OSM (Geoclimate) | Proportion |
| <i>Urban morphology</i> | Number of urban patches | OSM (Geoclimate) | Counts |
| | Mean surface of the urban patches | OSM (Geoclimate) | Positive continuous |
| | Proportion of high buildings | OSM (Geoclimate) | Proportion |
| | Proportion of mid-rise buildings | OSM (Geoclimate) | Proportion |
| | Proportion of compact buildings | OSM (Geoclimate) | Proportion |
| | Proportion of non-compact buildings | OSM (Geoclimate) | Proportion |
| <i>Human activity</i> | Population in the municipality | INSEE (municipality) | Counts |

MULTIVARIATE ANALYSIS



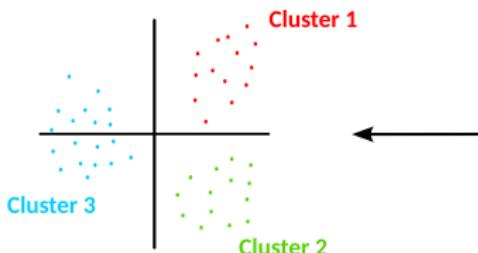
Factorial Analysis of Mixed Data

1/ Compute the first eigen-value for each group of variable

2/ Perform PCA on the weighted table

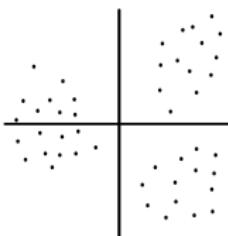
$$\left[\frac{X_1}{\sqrt{\lambda_1}}; \dots; \frac{X_j}{\sqrt{\lambda_j}}; \dots; \frac{X_J}{\sqrt{\lambda_J}} \right]$$

K-means clustering

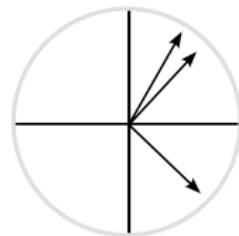


Project variables and individuals in their respective reduced-dimensional spaces

Projection of individuals



Projection of variables



VARIABLE PREPROCESSING

Counts and positive continuous variables

⇒ log-transformation: $y = \log(x + 1)$

VARIABLE PREPROCESSING

Counts and positive continuous variables

➡ log-transformation: $y = \log(x + 1)$

Proportions

➡ compositional data analysis

VARIABLE PREPROCESSING

Counts and positive continuous variables

- ➡ log-transformation: $y = \log(x + 1)$

Proportions

- ➡ compositional data analysis

Huge literature on the topic:

[Aitchison, 1982] The statistical analysis of compositional data. *Journal of the Royal Statistical Society. Series B.*

[Aitchison and J. Egozcue, 2005] Compositional Data Analysis: Where Are We and Where Should We Be Heading? *Mathematical Geology.*

[Greenacre, 2021] Compositional data analysis. *Annual Review of Statistics and its Application.*

→ Solution adopted by the statistical community: **logratio transformations**.

VARIABLE PREPROCESSING

Counts and positive continuous variables

⇒ log-transformation: $y = \log(x + 1)$

Proportions

⇒ compositional data analysis

Huge literature on the topic:

[Aitchison, 1982] The statistical analysis of compositional data. *Journal of the Royal Statistical Society. Series B.*

[Aitchison and J. Egozcue, 2005] Compositional Data Analysis: Where Are We and Where Should We Be Heading? *Mathematical Geology*.

[Greenacre, 2021] Compositional data analysis. *Annual Review of Statistics and its Application*.

→ Solution adopted by the statistical community: **logratio transformations**.

A very common one is the **centered logratio transformation (CLR)**:

$$\text{CLR}(j) = \log \left(\frac{x_j}{\left(\prod_{j'} x_{j'} \right)^{1/J}} \right) = \log(x_j) - \frac{1}{J} \sum_{j'} \log(x_{j'}) \quad j = 1, \dots, J.$$

VARIABLE PREPROCESSING

The CLR is not suited for zero-inflated compositional data.

VARIABLE PREPROCESSING

The CLR is not suited for **zero-inflated compositional data**.

A standard way to solve the issue → replace zeroes by small values.

VARIABLE PREPROCESSING

The CLR is not suited for **zero-inflated compositional data**.

A standard way to solve the issue → replace zeroes by small values.

Following [Martín-Fernández et al., 2003], a multiplicative replacement strategy gives:

$$r_j = \begin{cases} \delta_j, & \text{if } x_j = 0, \\ \left(1 - \frac{\sum_{k|x_k=0}\delta_k}{c}\right)x_j, & \text{if } x_j > 0, \end{cases}$$

where

- δ_j is the imputed value on the part x_j
- c is the constant of the sum constraint ($c = 1$ if x_j are proportions, $c = 100$ if x_j are percentage).

VARIABLE PREPROCESSING

The CLR is not suited for **zero-inflated compositional data**.

A standard way to solve the issue → replace zeroes by small values.

Following [Martín-Fernández et al., 2003], a multiplicative replacement strategy gives:

$$r_j = \begin{cases} \delta_j, & \text{if } x_j = 0, \\ \left(1 - \frac{\sum_{k|x_k=0}\delta_k}{c}\right)x_j, & \text{if } x_j > 0, \end{cases}$$

where

- δ_j is the imputed value on the part x_j
- c is the constant of the sum constraint ($c = 1$ if x_j are proportions, $c = 100$ if x_j are percentage).

A reasonable choice for the imputed value is to take $\delta_j = \frac{2}{3} \min(x^*)$ where x^* is the positive part of the composition [Greenacre, 2021].

VARIABLE PREPROCESSING

The CLR is not suited for **zero-inflated compositional data**.

A standard way to solve the issue → replace zeroes by small values.

Following [Martín-Fernández et al., 2003], a multiplicative replacement strategy gives:

$$r_j = \begin{cases} \delta_j, & \text{if } x_j = 0, \\ \left(1 - \frac{\sum_{k|x_k=0}\delta_k}{c}\right)x_j, & \text{if } x_j > 0, \end{cases}$$

where

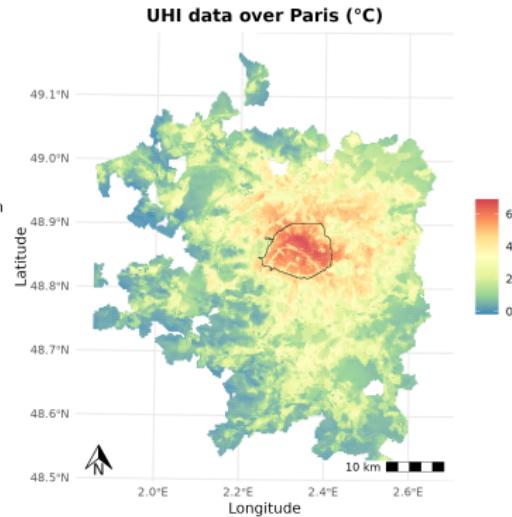
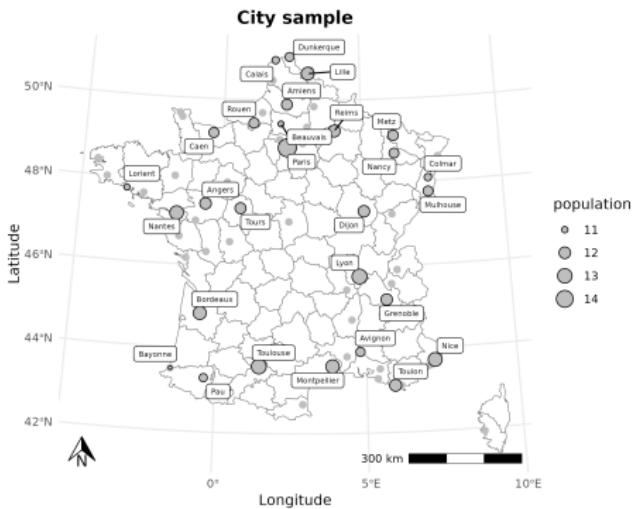
- δ_j is the imputed value on the part x_j
- c is the constant of the sum constraint ($c = 1$ if x_j are proportions, $c = 100$ if x_j are percentage).

A reasonable choice for the imputed value is to take $\delta_j = \frac{2}{3} \min(\mathbf{x}^*)$ where \mathbf{x}^* is the positive part of the composition [Greenacre, 2021].

► This zero-replacement strategy implies that zeroes are not considered to be **essential zeros**, but rather values that are too small to be detected.

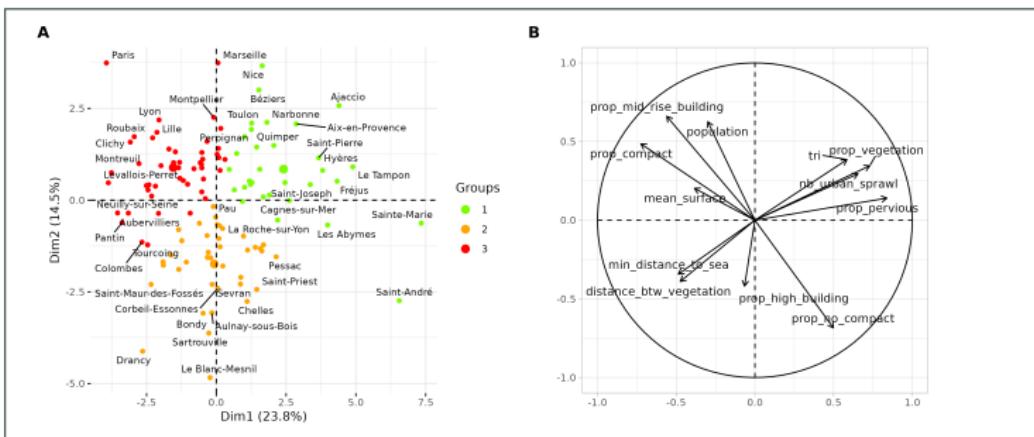
CITY SAMPLE

- Municipalities over 50 000 inhabitants
 - ➡ 120 municipalities (1 individual = 1 municipality)
- UHI model outputs from TEB/SURFEX model (MaPuce project)
 - ➡ 42 urban areas



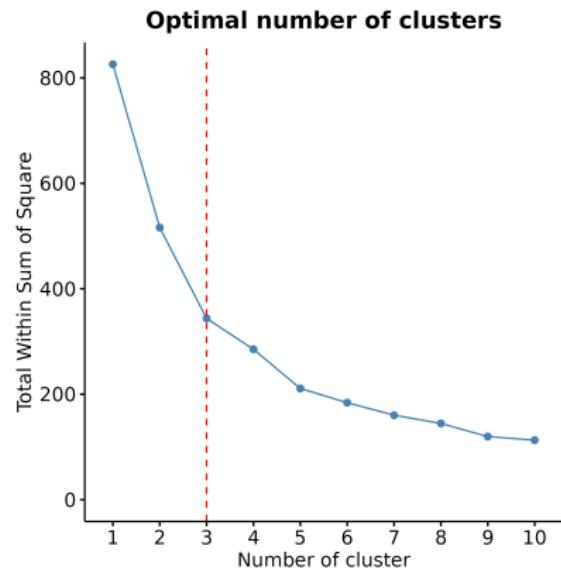
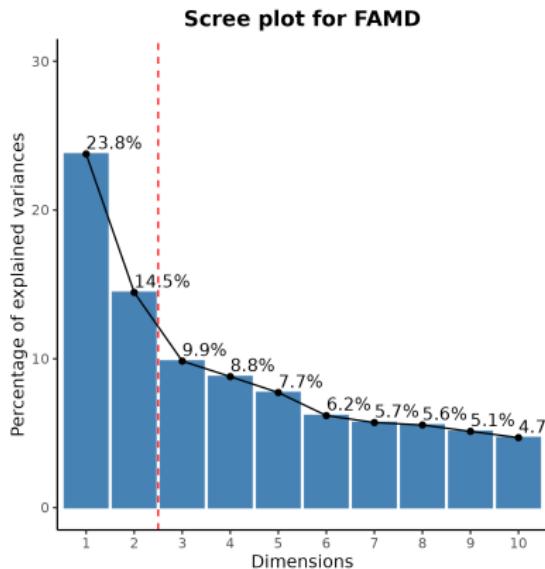
RESULTS

FAMD and clustering results



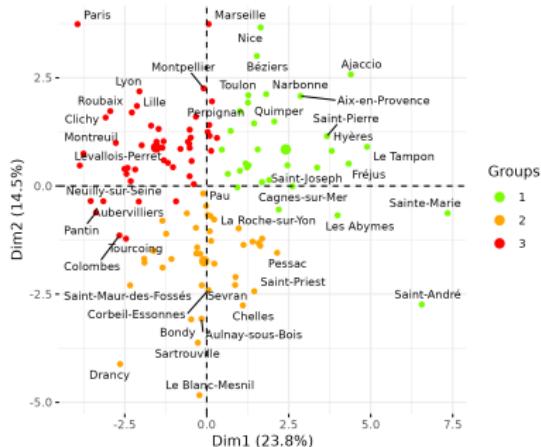
RESULTS

1. FAMD → filter two first dimensions.
2. Conduct k-means based on the two first dimensions
→ three clusters.

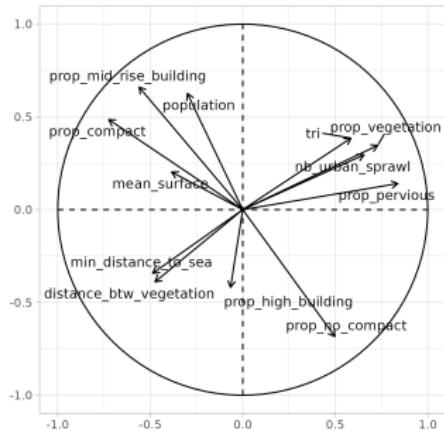


RESULTS

A



B



Structuring variables :

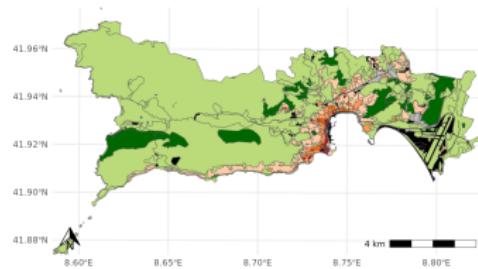
- ➡ Density and height of the building.
- ➡ "Green" variables, number of urban patches, topography.

Two main axis of variability:

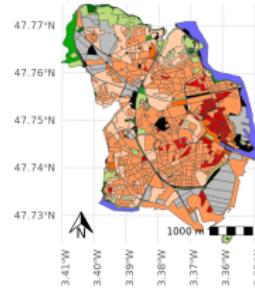
- 1/ Vegetalization and cooling variables.
- 2/ Compact versus non-compact.

RESULTS

Cluster 1
Ajaccio



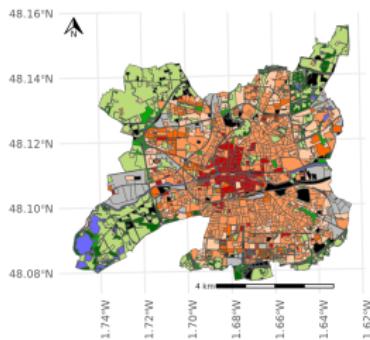
Cluster 2
Lorient



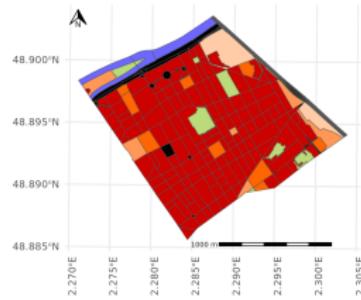
LCZ

- LCZ 1: Compact high-rise
- LCZ 2: Compact mid-rise
- LCZ 3: Compact low-rise
- LCZ 4: Open high-rise
- LCZ 5: Open mid-rise
- LCZ 6: Open low-rise
- LCZ 7: Lightweight low-rise
- LCZ 8: Large low-rise
- LCZ 9: Sparsely built
- LCZ 10: Heavy industry
- LCZ A: Dense tress
- LCZ B: Scattered trees
- LCZ C: Bush, Scrub
- LCZ D: Low plants
- LCZ E: Bare rock or Paved
- LCZ F: Bare soil or Sand
- LCZ G: Water

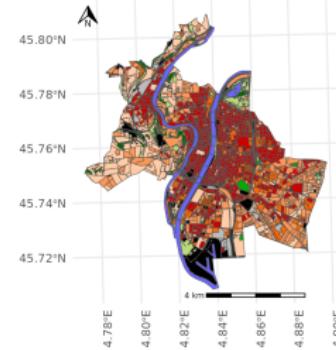
Cluster 3
Rennes



Cluster 3
Levallois-Perret

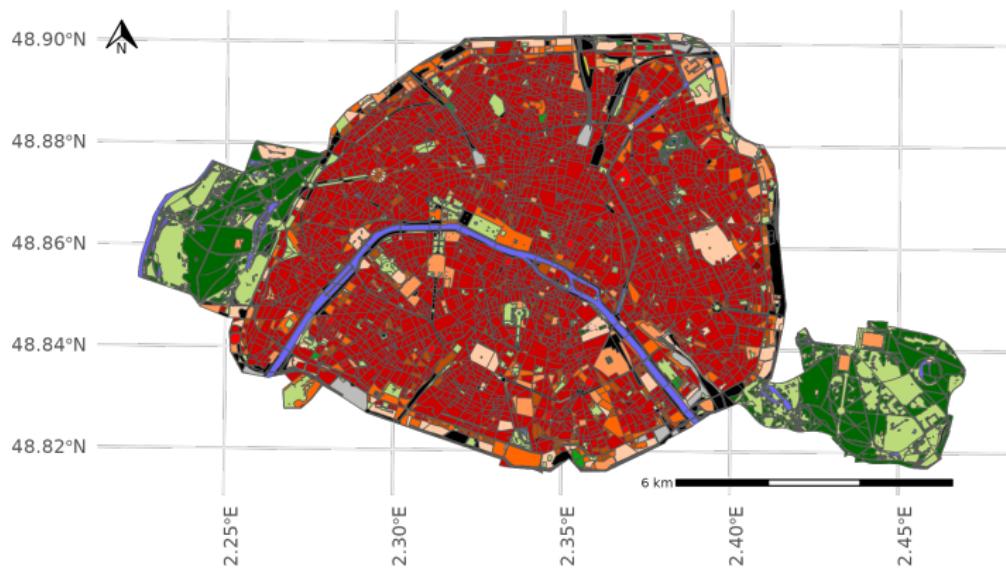


Cluster 3
Lyon

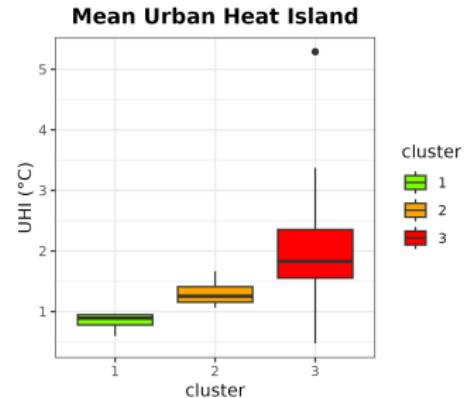
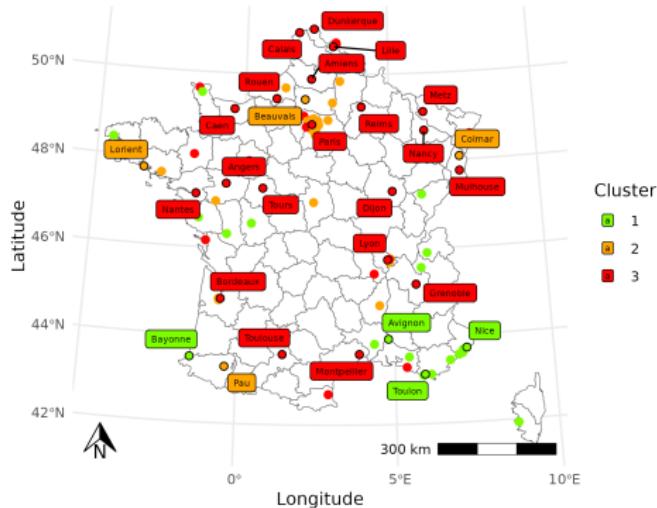


RESULTS

Paris = outlier



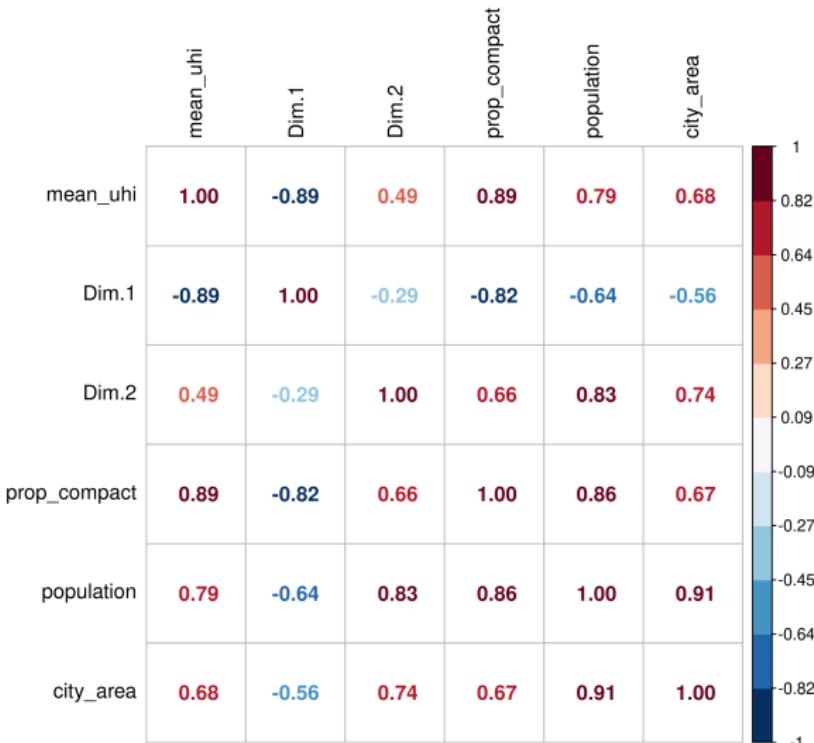
RESULTS



- **Cluster 1:** coastal or mountainous, high permeability, several urban patches.
- **Cluster 2:** single urban sprawl, low permeability, non-compact LCZ.
- **Cluster 3:** dense urban fabric, many urban barriers, minimal permeability.

RESULTS

Correlation between the mean urban heat island and some key variables of the analysis.



DISCUSSION

We identified (1) municipality profiles with regards to overheating factors and (2) the main variables that drive these profiles.

DISCUSSION

We identified (1) municipality profiles with regards to overheating factors and (2) the main variables that drive these profiles.

The intent of this work is to propose a **methodological framework** to identify such city profiles.

Steps of the methodological framework

1. For each overheating factor, identify geographic indicators that describe these factors.
2. Select the urban territories of interest (administrative units, geographical extent) and compute the geographic indicators.
3. Identify profiles from the data through statistical analysis.
4. Confront these to some physical measure of overheating.

DISCUSSION

We identified (1) municipality profiles with regards to overheating factors and (2) the main variables that drive these profiles.

The intent of this work is to propose a **methodological framework** to identify such city profiles.

Steps of the methodological framework

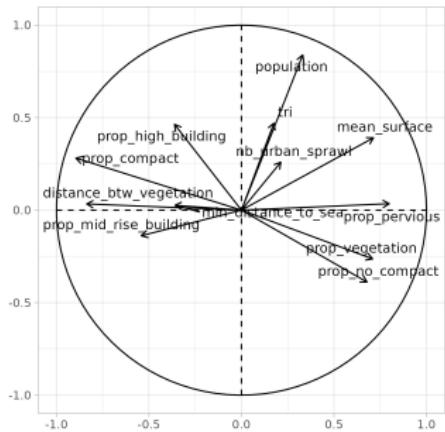
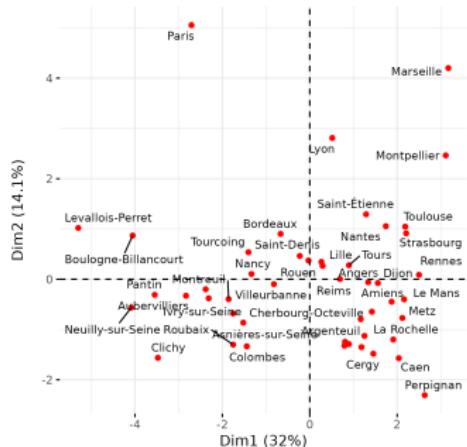
1. For each overheating factor, identify geographic indicators that describe these factors.
2. Select the urban territories of interest (administrative units, geographical extent) and compute the geographic indicators.
3. Identify profiles from the data through statistical analysis.
4. Confront these to some physical measure of overheating.

From our data, we are able to build interpretable city profiles relative to overheating factors, but this is not a definitive classification.

Many choices in this work are open to debate and we will discuss these now.

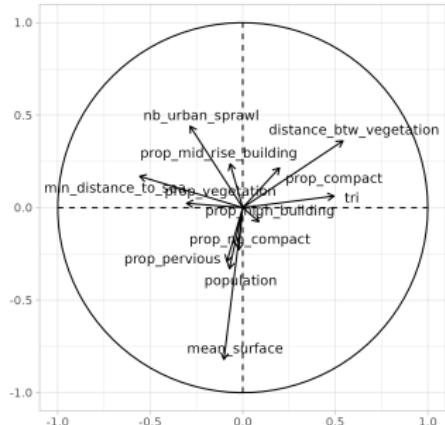
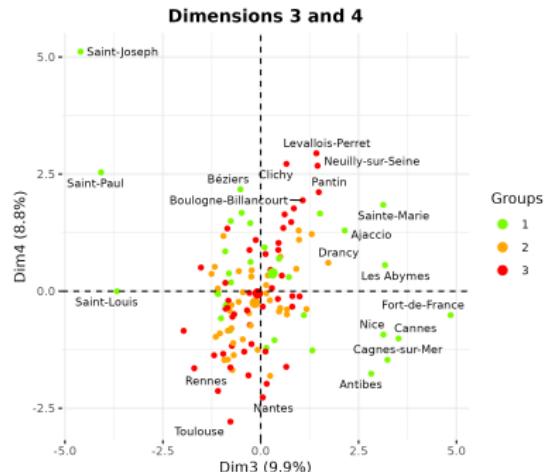
EXPLORING THE CLUSTER OF HEAT-SENSITIVE CITIES

Dimensions 1 and 2



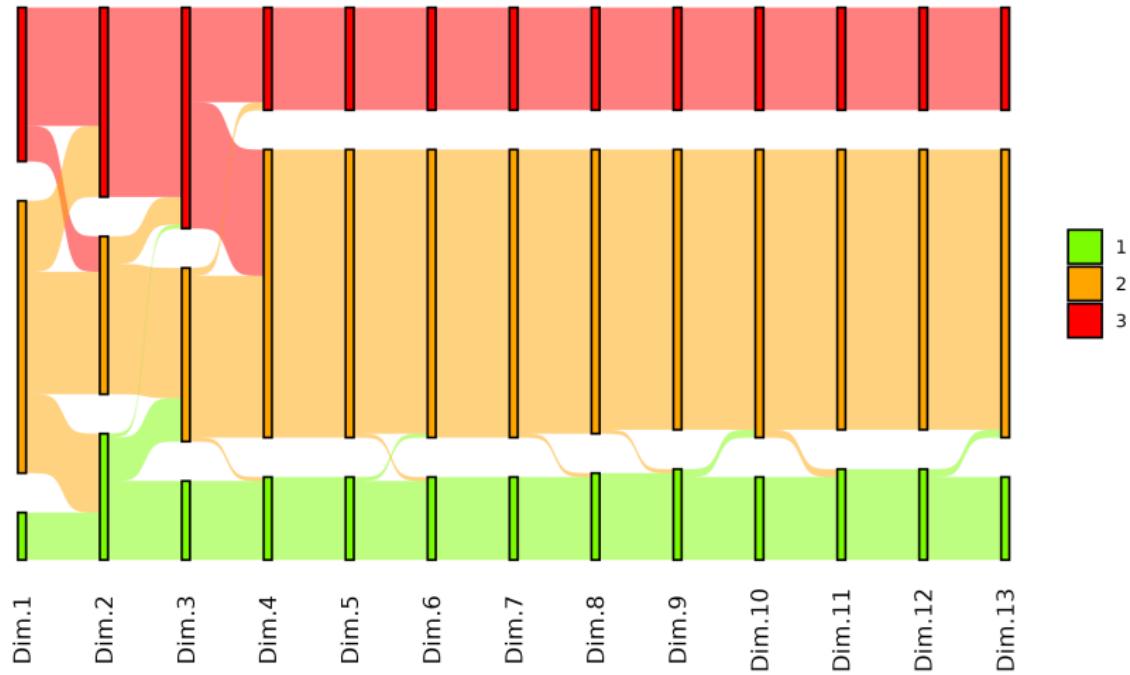
- known relationship: non-compact versus compact
- and some new one: 'population' distinguishes from 'compact'

EXPLORING THE NEXT DIMENSIONS OF THE FAMD

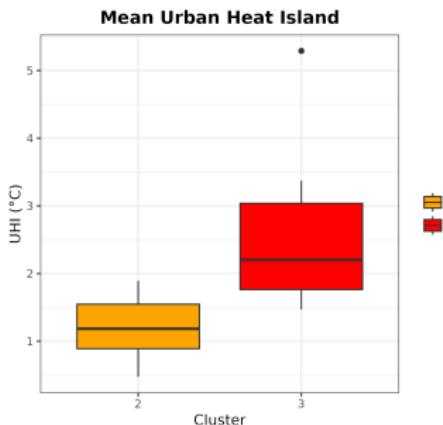
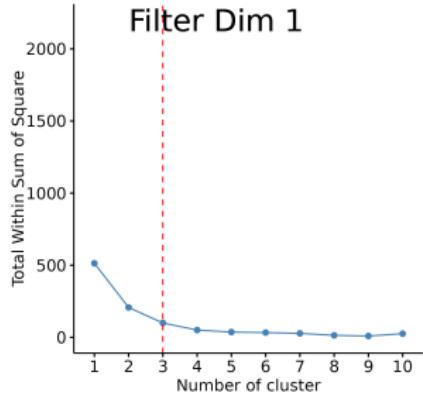
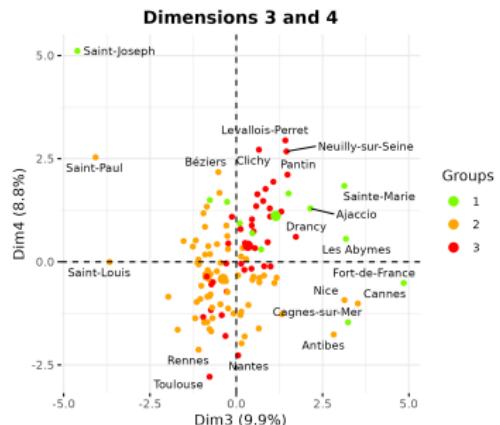
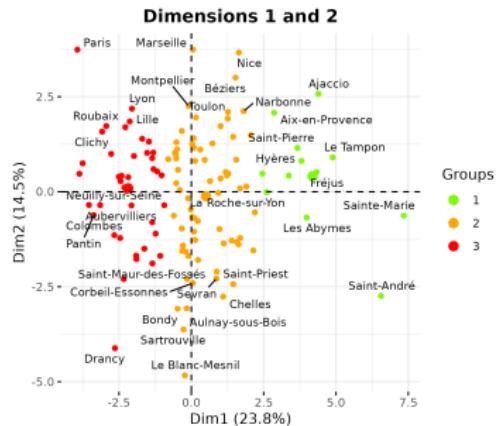


- 3rd dimension: topography and distance to the sea
- 4th dimension: surface of the urban patches

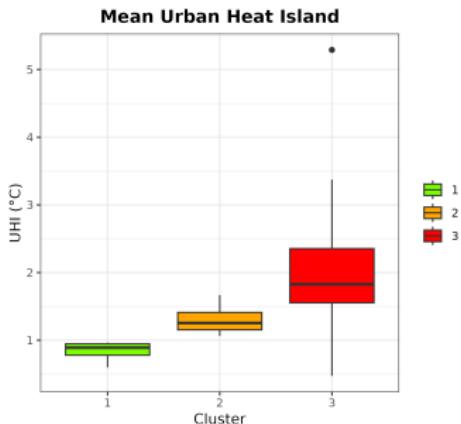
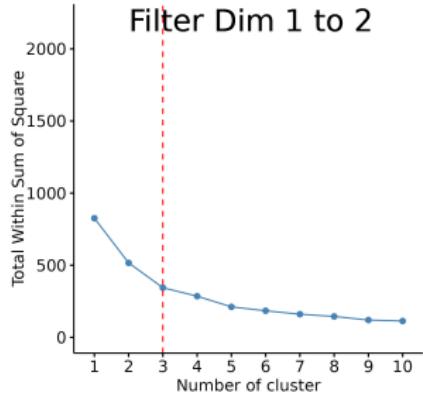
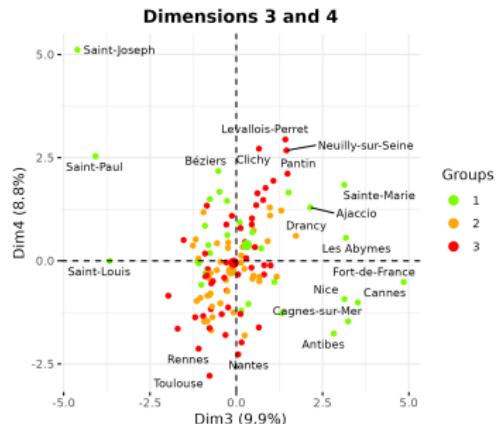
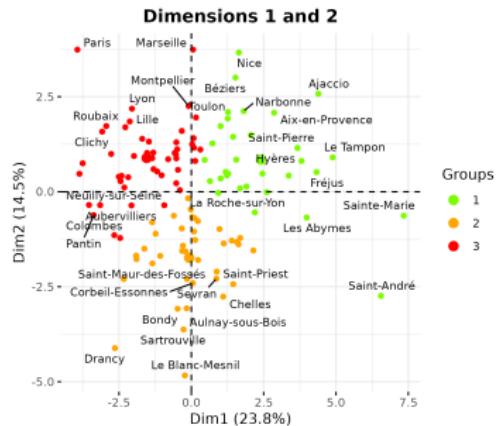
EFFECT OF THE NUMBER OF AFMD DIMENSIONS IN THE CLUSTERING



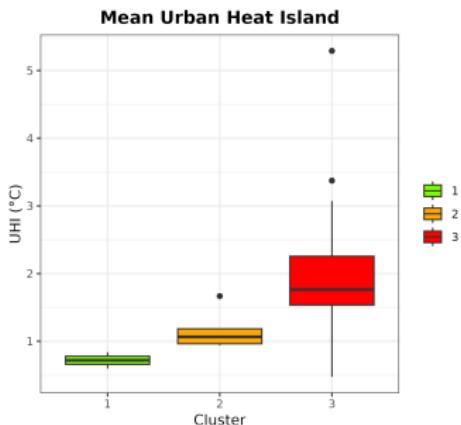
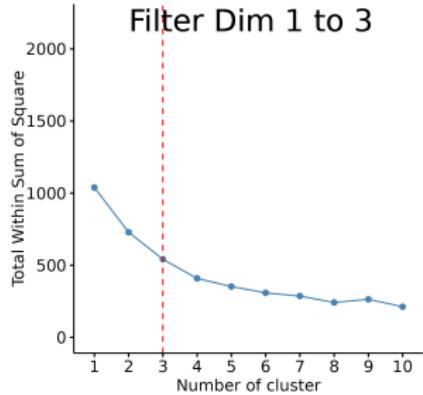
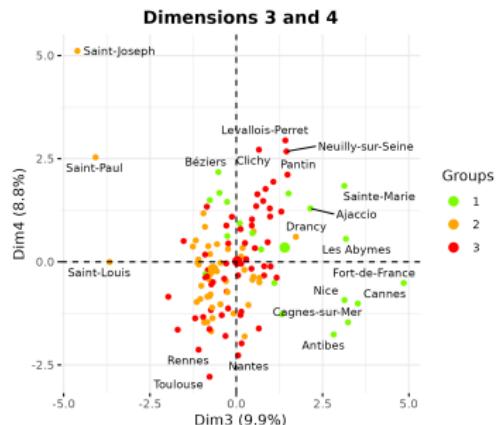
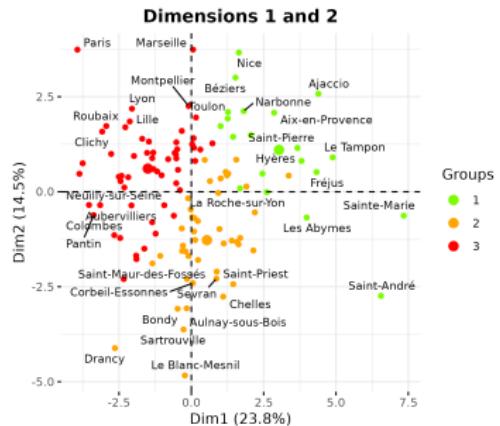
EFFECT OF THE NUMBER OF AFMD DIMENSIONS IN THE CLUSTERING



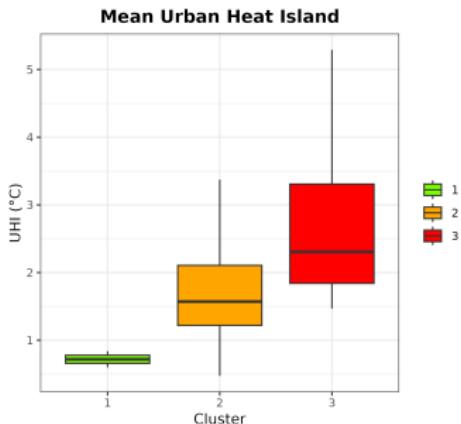
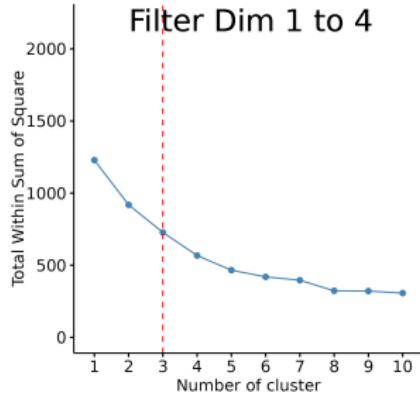
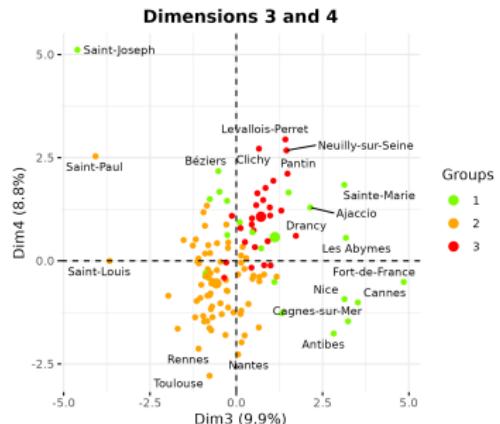
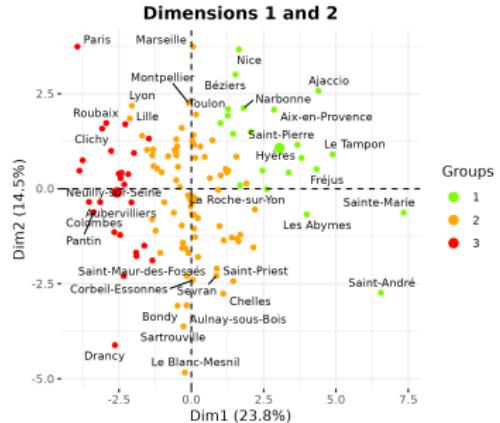
EFFECT OF THE NUMBER OF AFMD DIMENSIONS IN THE CLUSTERING



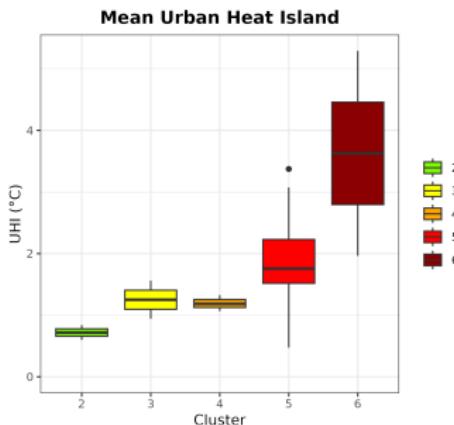
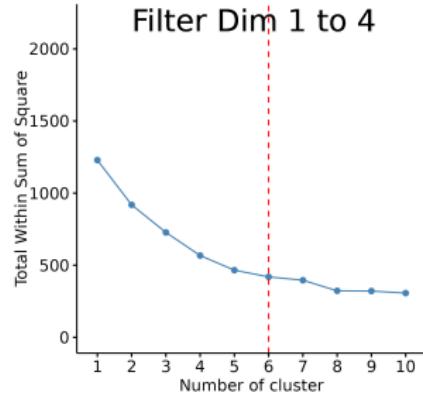
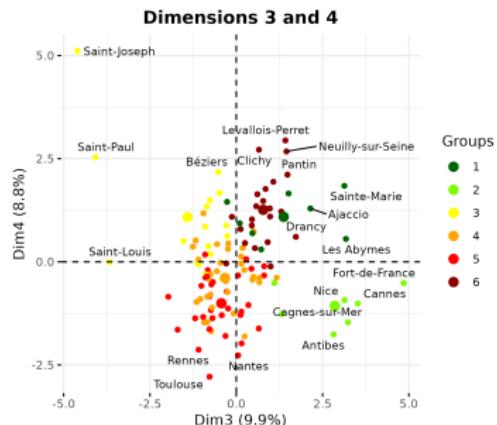
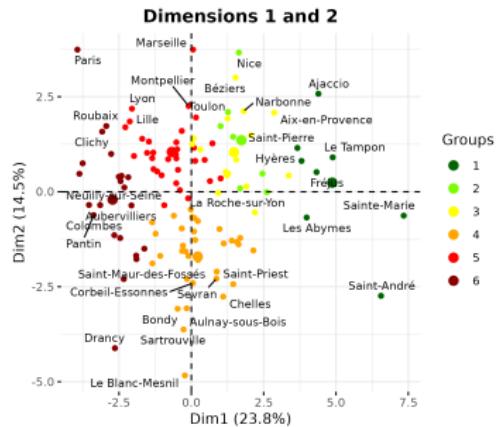
EFFECT OF THE NUMBER OF AFMD DIMENSIONS IN THE CLUSTERING



EFFECT OF THE NUMBER OF AFMD DIMENSIONS IN THE CLUSTERING



PROPOSE A FINER CLASSIFICATION



EXTENDING THE ANALYSIS

Geographically

Consider a wider area (Europe, Africa, Asia, America)

Administratively

Consider urban units rather than municipalities

By including the within structure of the city

The within-structure of city can affect overheating.

→ build geographical indicators representing these structures

➡ More robust and exhaustive classification

A PROBABILISTIC PCA FOR HANDLING ZERO VALUES

The zero-replacement strategy for compositions relies on some arbitrary choices.

The main advances to account for essential zeroes in compositions come from microbial research where microbiome data are considered as compositions [Zeng et al., 2023].

A PROBABILISTIC PCA FOR HANDLING ZERO VALUES

The zero-replacement strategy for compositions relies on some arbitrary choices.

The main advances to account for essential zeroes in compositions come from microbial research where microbiome data are considered as compositions [Zeng et al., 2023].

Consider a latent variable model where \mathbf{Y} is the matrix of variables and \mathbf{Z} are the latent variable. $\boldsymbol{\theta}$ and σ are the parameters.

$$\mathbf{Y} | \mathbf{Z}, \boldsymbol{\theta}, \sigma \sim \mathcal{L}(\mathbf{Z}, \sigma^2 \mathbf{I}),$$

$$\mathbf{Z} | \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

$$\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}'$$

where \mathbf{B} is low rank ($n_{col}(\mathbf{B}) < n_{col}(\mathbf{Y})$).

A PROBABILISTIC PCA FOR HANDLING ZERO VALUES

The zero-replacement strategy for compositions relies on some arbitrary choices.

The main advances to account for essential zeroes in compositions come from microbial research where microbiome data are considered as compositions [Zeng et al., 2023].

Consider a latent variable model where \mathbf{Y} is the matrix of variables and \mathbf{Z} are the latent variable. $\boldsymbol{\theta}$ and σ are the parameters.

$$\mathbf{Y} | \mathbf{Z}, \boldsymbol{\theta}, \sigma \sim \mathcal{L}(\mathbf{Z}, \sigma^2 \mathbf{I}),$$

$$\mathbf{Z} | \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

$$\boldsymbol{\Sigma} = \mathbf{B}\mathbf{B}'$$

where \mathbf{B} is low rank ($n_{col}(\mathbf{B}) < n_{col}(\mathbf{Y})$).

Do we model zeroes in $\mathcal{L}()$ or in \mathbf{Z} ?

→ in $\mathcal{L}()$: several zero-inflated distributions such as zero-inflated dirichlet, beta, logistic-normal.

→ in \mathbf{Z} : modeling zeroes and positive values through delta models.

TAKE HOME MESSAGE

Main points:

- **Very few classifications** of cities based on their sensitivity to overheating.
- Here, **classification based on 120 French municipalities** using 13 indicators related to climate, land use, morphology, and human activity.
- Factorial Analysis of Mixed Data (FAMD) \times k-means clustering \rightarrow **3 urban profiles** with differing levels of vegetation, density, and permeability.
- These clusters show **strong relation with UHI model outputs**.

Limits and perspectives:

- Extend to other countries/continents
- Look at a finer scale to the inner structure of the city
- Adapt the methodology to better account for heterogeneous data types
 \rightarrow probabilistic PCA

Thank you for your attention!

Happy to take questions? 😊

Internship application to PEPR Maths-Vives - Any candidate ?

**SCIENCE
DES DONNÉES
HISTOIRE
& TERRITOIRES**

10 ÈMES RENCONTRES DE STATISTIQUE



27 & 28

NOVEMBRE 2025

Amphithéâtre Yves Coppens
Faculté Sciences & Sciences
de l'Ingénieur
Université Bretagne Sud
Campus de Tohannic - VANNES

REFERENCES I

-  Aitchison, J. (1982).
The statistical analysis of compositional data.
Journal of the Royal Statistical Society: Series B (Methodological), 44(2):139–160.
-  Aitchison, J. and J. Egozcue, J. (2005).
Compositional data analysis: where are we and where should we be heading?
Mathematical Geology, 37(7):829–850.
-  Bocher, E., Bernard, J., Wiederhold, E. L. S., Leconte, F., Petit, G., Palominos, S., and Noûs, C. (2021).
Geoclimate: a geospatial processing toolbox for environmental and climate studies.
Journal of Open Source Software, 6(65):3541.

REFERENCES II

-  Cerema (2024).
Local climate zones dashboard.
https://cartagene.cerema.fr/portal/apps/dashboards/08066acd23974111be1584a5761fd6b9.
Accessed: 2025-11-18.
-  Ching, J., Mills, G., Bechtel, B., See, L., Feddema, J., Wang, X., Ren, C., Brousse, O., Martilli, A., Neophytou, M., et al. (2018).
Wudapt: An urban weather, climate, and environmental modeling infrastructure for the anthropocene.
Bulletin of the American Meteorological Society, 99(9):1907–1924.
-  Greenacre, M. (2021).
Compositional data analysis.
Annual Review of Statistics and its Application, 8(1):271–299.

REFERENCES III

-  Huang, W. T. K., Masselot, P., Bou-Zeid, E., Faticchi, S., Paschalis, A., Sun, T., Gasparrini, A., and Manoli, G. (2023).
Economic valuation of temperature-related mortality attributed to urban heat islands in european cities.
Nature communications, 14(1):7438.
-  Husson, F., Lê, S., and Pagès, J. (2011).
Exploratory multivariate analysis by example using R, volume 15.
CRC press Boca Raton.
-  Martín-Fernández, J. A., Barceló-Vidal, C., and Pawlowsky-Glahn, V. (2003).
Dealing with zeros and missing values in compositional data sets using nonparametric imputation.
Mathematical Geology, 35(3):253–278.

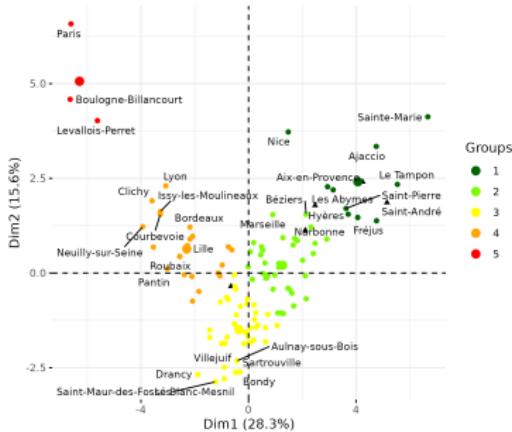
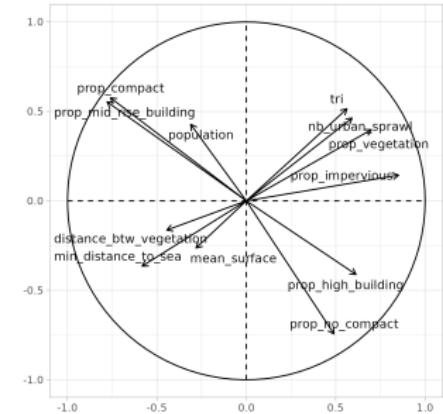
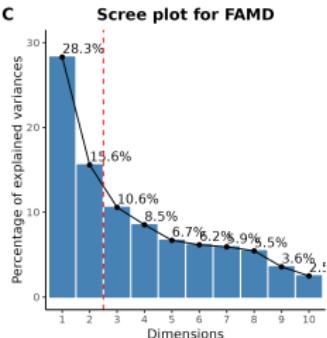
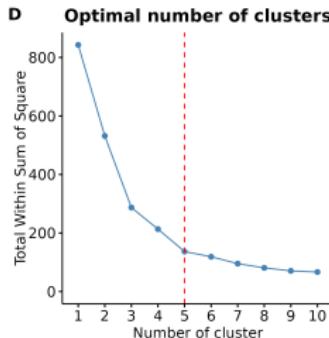
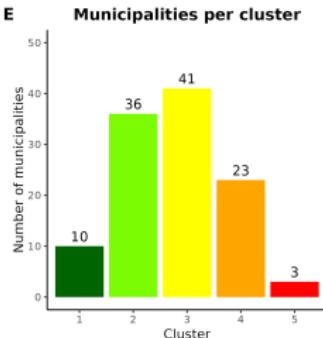
REFERENCES IV

-  Masson, V. (2000).
A physically-based scheme for the urban energy budget in atmospheric models.
Boundary-layer meteorology, 94:357–397.
-  Vargas-Munoz, J. E., Srivastava, S., Tuia, D., and Falcao, A. X. (2020).
Openstreetmap: Challenges and opportunities in machine learning and remote sensing.
IEEE Geoscience and Remote Sensing Magazine, 9(1):184–199.
-  Wang, C., Wang, Z.-H., and Li, Q. (2020).
Emergence of urban clustering among us cities under environmental stressors.
Sustainable Cities and Society, 63:102481.

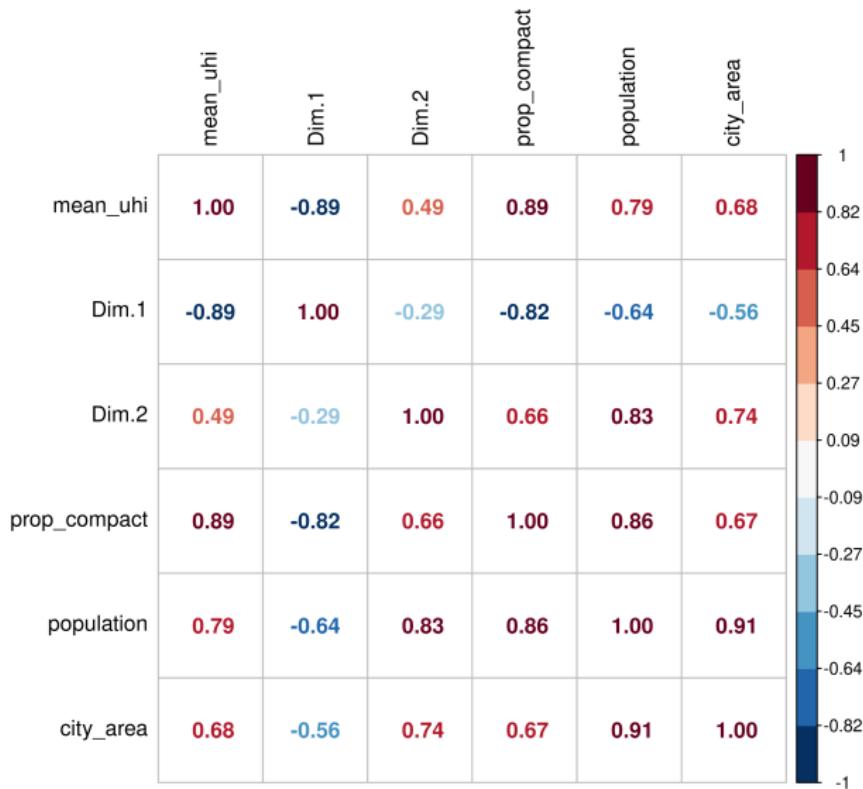
REFERENCES V

-  Zeng, Y., Pang, D., Zhao, H., and Wang, T. (2023).
A zero-inflated logistic normal multinomial model for extracting microbial compositions.
Journal of the American Statistical Association,
118(544):2356–2369.

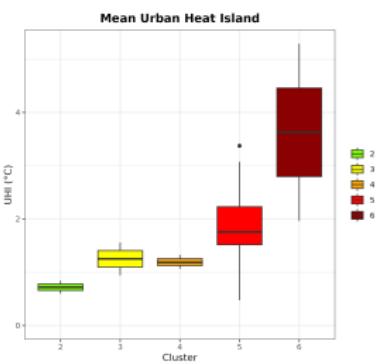
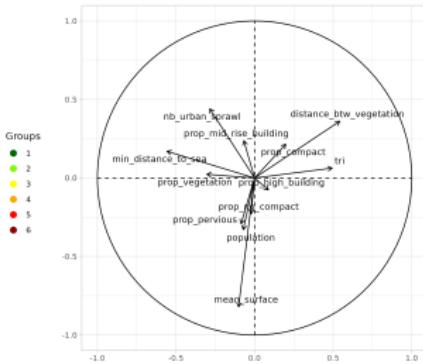
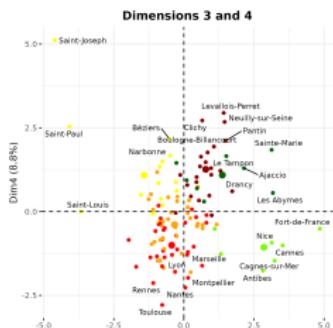
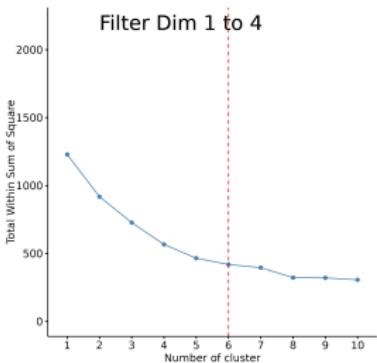
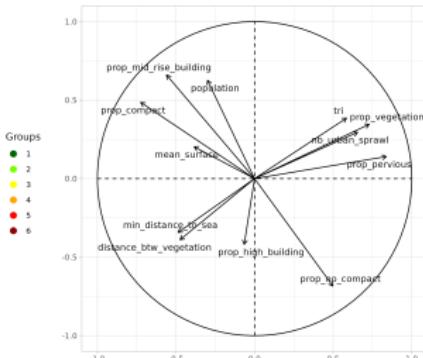
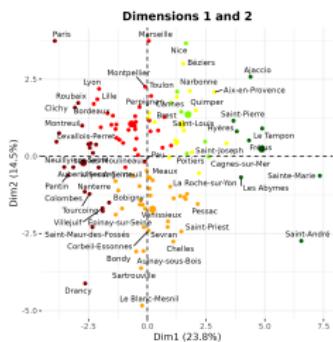
APPENDICES

A**B****C****D****E**

APPENDICES

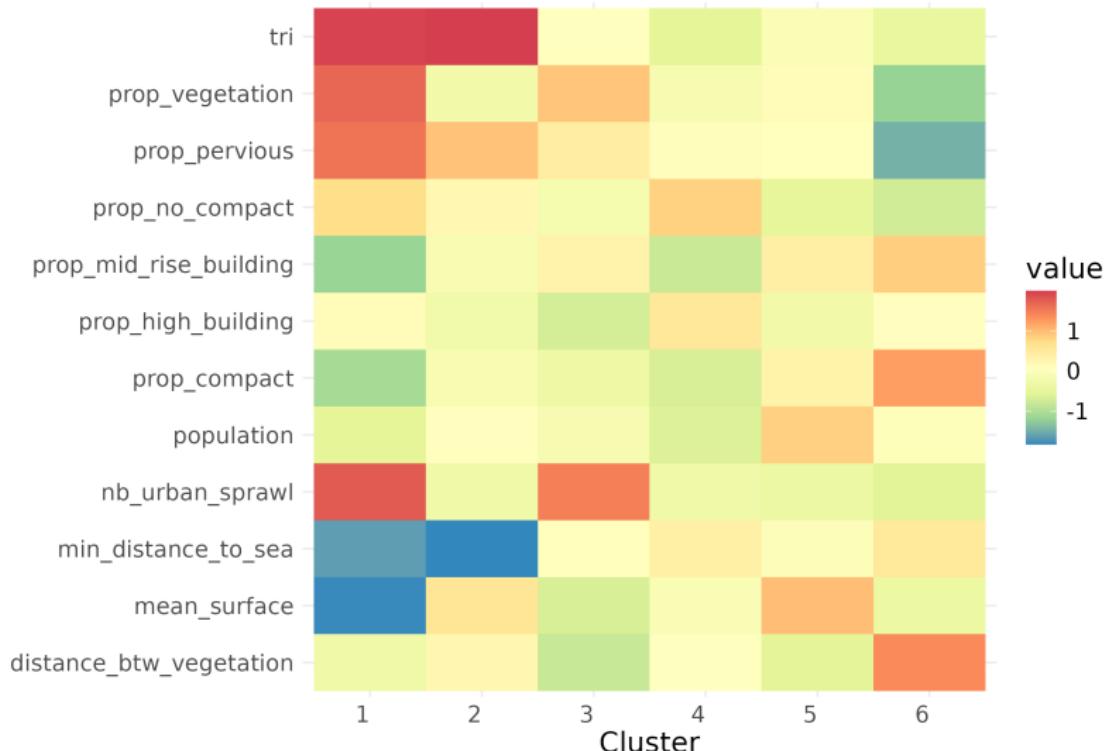


APPENDICES

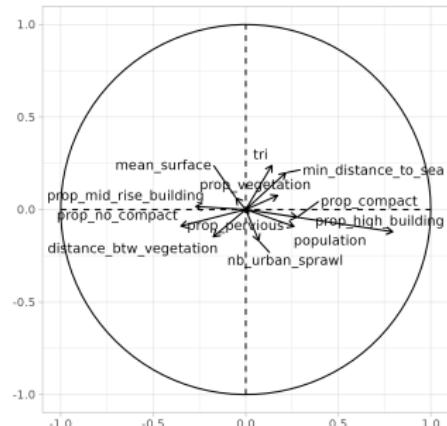
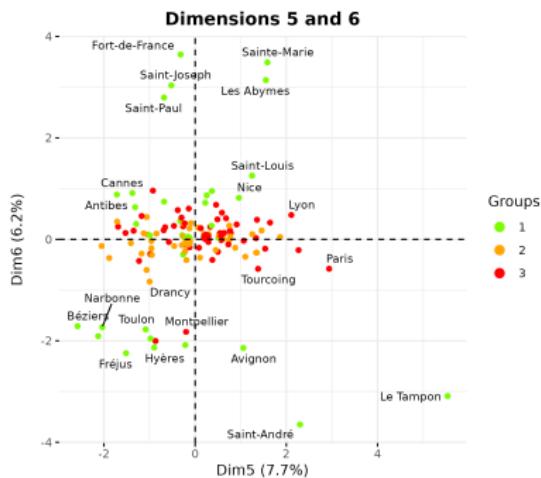


APPENDICES

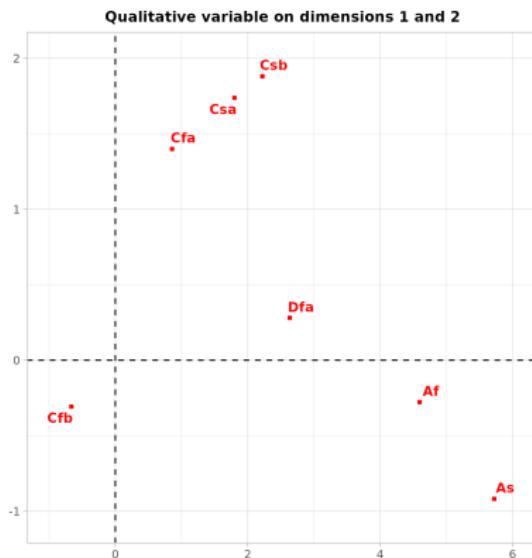
Standardized variable value for each cluster



APPENDICES



APPENDICES



Köppen-Geiger Classes Shown in the Plot

