



# Positive Matrix Factorization and GIS approach to perform data mining on groundwater and surface water quality dataset.

Zanotti, C.<sup>1</sup>, Rotiroti, M.<sup>1</sup>, Fumagalli, L.<sup>1</sup>, Stefania, G.A.<sup>1</sup>, Canonaco, F.<sup>2</sup>, Stefenelli, G.<sup>2</sup>, Prévôt, A.S.H.<sup>2</sup>, Leoni, B.<sup>1</sup>, Bonomi, T.<sup>1</sup>.

<sup>1</sup> Department of Earth and Environmental Sciences, University of Milano-Bicocca, Piazza della Scienza, 1, 20126 Milano, Italy. Email: [c.zanotti@campus.unimib.it](mailto:c.zanotti@campus.unimib.it)

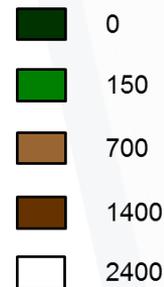
<sup>2</sup> Laboratory of Atmospheric Chemistry, Paul Scherrer Institute, 5232 Villigen-PSI, Switzerland

## Main aims:

- To identify the main hydrochemical features of either groundwater and surface water (lake, river and springs) and, possibly, the processes that influence (or govern) them.
- To understand the relations (if any) between the chemistry of groundwater and surface water.

## Methodology:

- Application of Positive Matrix Factorization (PMF)

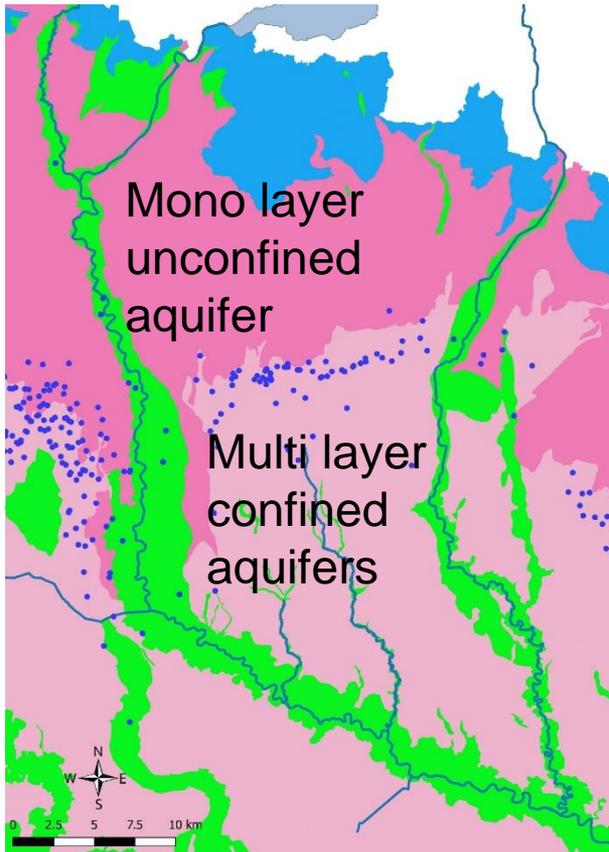


**Oglio river catchment area:** 6360 km<sup>2</sup>

**Study area:** 1960 km<sup>2</sup>

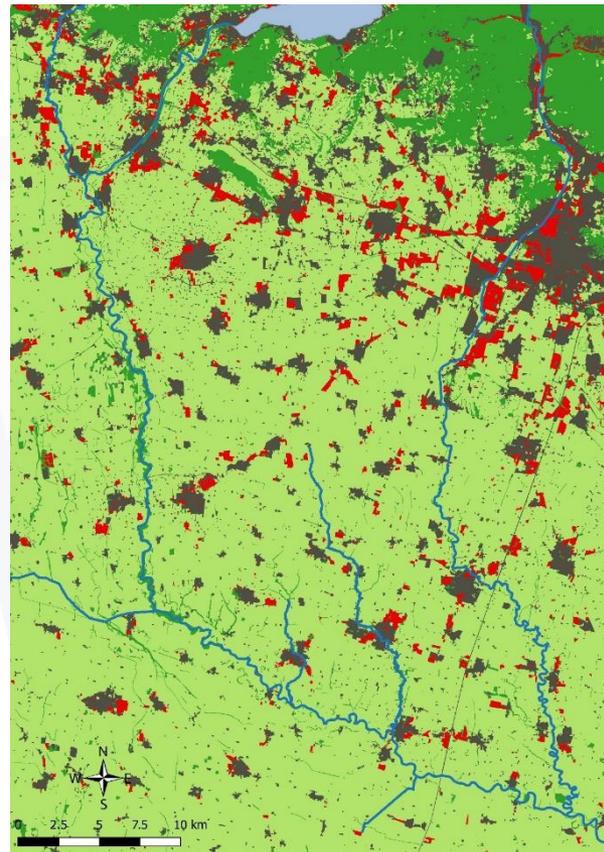
# 1 – STUDY AREA

## Geology



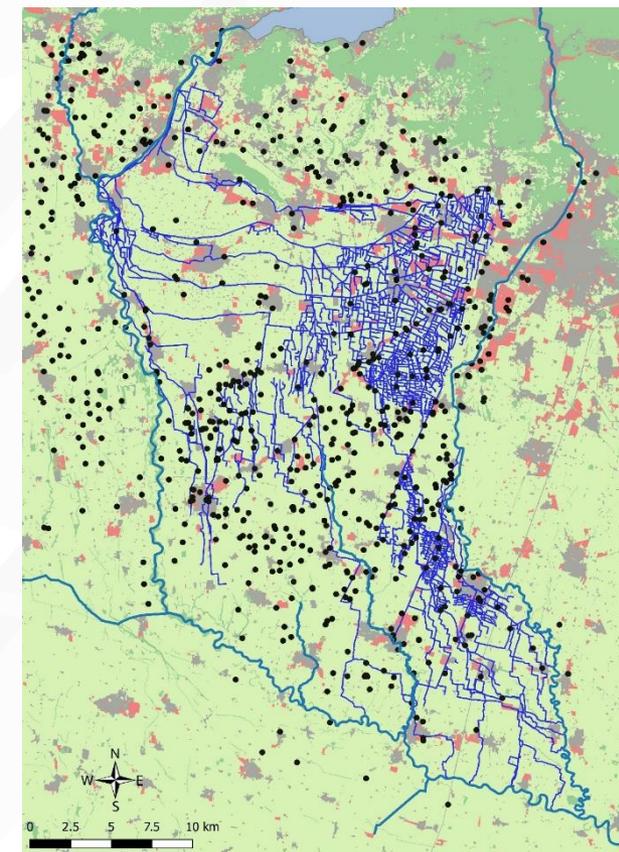
- Moraine
- Higher Plain
- Low Plain
- River Valley
- springs

## Landuse



- Agriculture
- Urban
- Industry
- Natural

## Irrigation sources



- Irrigation channels
- Irrigation wells

## 4 field surveys:

- February 2016
- June 2016
- September 2016
- March 2017.

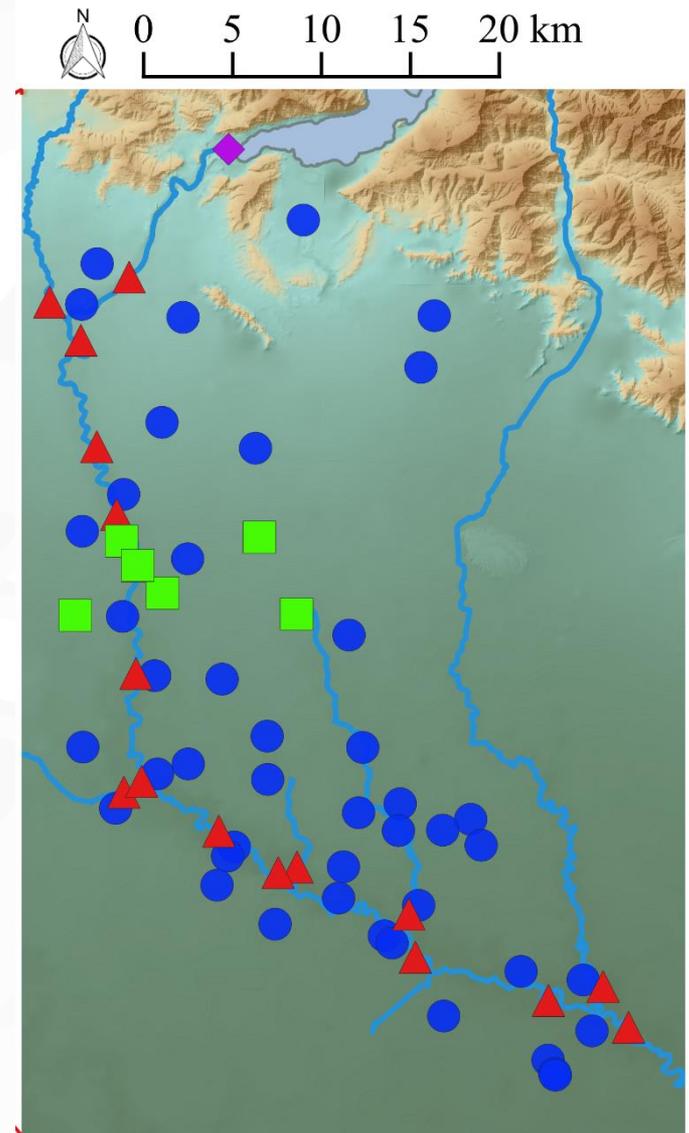
Measured in water samples from:

- ● 44 wells
- ◆ 1 lake station
- ▲ 17 river stations
- ■ 6 springs

Chemical variables (14):

- Dissolved O<sub>2</sub>
- P<sub>tot</sub>
- Major ions (HCO<sub>3</sub>, Ca, Mg, Na, K, Cl, SO<sub>4</sub>, NO<sub>3</sub>, NH<sub>4</sub>)
- Trace elements (As, Fe, Mn)

For the validation: Water Isotopes





Problem related to the reduction the original data-set into one of lower dimensions to detect “hidden” information and explain the variability of the measured variables.

Usually this kind of problem is solved with Factor Analysis, but it presents some limitations:

## Factor Analysis (FA)



- Well known and standardized
- Widely used in hydrology



- Factors grouping variables with both positive and negative correlation
- It is sometimes difficult to interpret and disseminate
- No way to treat data uncertainty

Positive Matrix Factorization (PMF) is a multivariate analysis aimed at source identification and apportionment, specifically designed to cope with environmental data and manage their uncertainty and distributions. Here it is tested the effectiveness of PMF as a tool to perform data mining and define hydrochemical features of groundwater and surface water and to understand their relationship.



## Positive Matrix Factorization (PMF)

- Positivity constraint
- Data uncertainty



## How does PMF work in hydrology?

## Positive Matrix Factorization (PMF)

PMF (Paatero and Tapper 1994) is a multivariate analysis in which the fundamental problem is to resolve the identities and contributions of different sources in a mixture.

$$x_{ij} = \sum_{k=1}^p g_{ik}f_{kj} + e_{ij} \quad i = 1 \dots m; j = 1 \dots n; k = 1 \dots p$$

➤ NO orthogonality constraint

➤ **2 Positivity constraints:**

- the composition of predicted source must be positive (a source cannot have negative percentage of an element);
- the predicted source contribution to each sample must all be positive (a source cannot contribute to a sample with a negative mass).

➤ The solution of the problem is found through a weighted least squares approach, minimizing:

$$Q = \sum_{i=1}^n \sum_{j=1}^m \left( \frac{x_{ij} - \sum_{k=1}^p g_{ik}f_{kj}}{\sigma_{ij}} \right)^2$$

**Data uncertainty**

Allows to work with missing data and data below detection limit



**Particularly suitable for environmental data**





## Informatic tool:

Atmos. Meas. Tech., 6, 3649–3661, 2013  
www.atmos-meas-tech.net/6/3649/2013/  
doi:10.5194/amt-6-3649-2013  
© Author(s) 2013. CC Attribution 3.0 License.



Atmospheric  
Measurement  
Techniques



## **SoFi, an IGOR-based interface for the efficient use of the generalized multilinear engine (ME-2) for the source apportionment: ME-2 application to aerosol mass spectrometer data**

**F. Canonaco, M. Crippa, J. G. Slowik, U. Baltensperger, and A. S. H. Prévôt**

Paul Scherrer Institute, Laboratory of Atmospheric Chemistry, 5232 Villigen PSI, Switzerland

*Correspondence to:* A. S. H. Prévôt (andre.prevot@psi.ch)



## Two major sources of uncertainty in PMF analysis:

- Rotational ambiguity (more than one solution can give the same Q)
- Random errors in the data



## To assess PMF uncertainty:

- 100 PMF runs over bootstrapped data
- At each run factor profiles rotates within predefined range

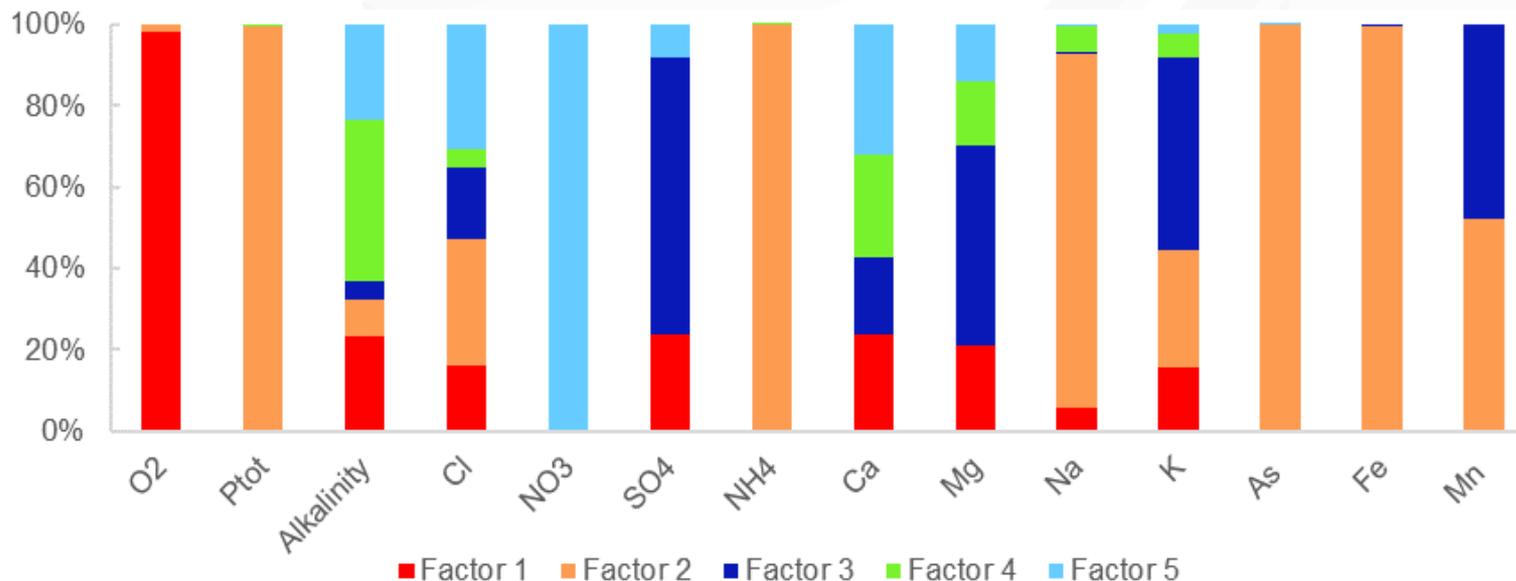
Results are expressed as average over the different runs and the uncertainty as the standard deviation.

At each run only 70% of the data are randomly resampled from the original dataset and used to solve the model  
This reduces the leverage of single data, especially extremes

# 1 – RESULTS – PMF factors profiles

## PMF generated 5 Factors

Each one of them was associated to a specific process or feature of the system, based on their factor profile or their spatial distribution



1- **Lake Iseo and Oglio river:** highly oxygenated water, with SO<sub>4</sub> and a small content of major ions such as HCO<sub>3</sub>, Cl, Ca, Mg, Na, and K

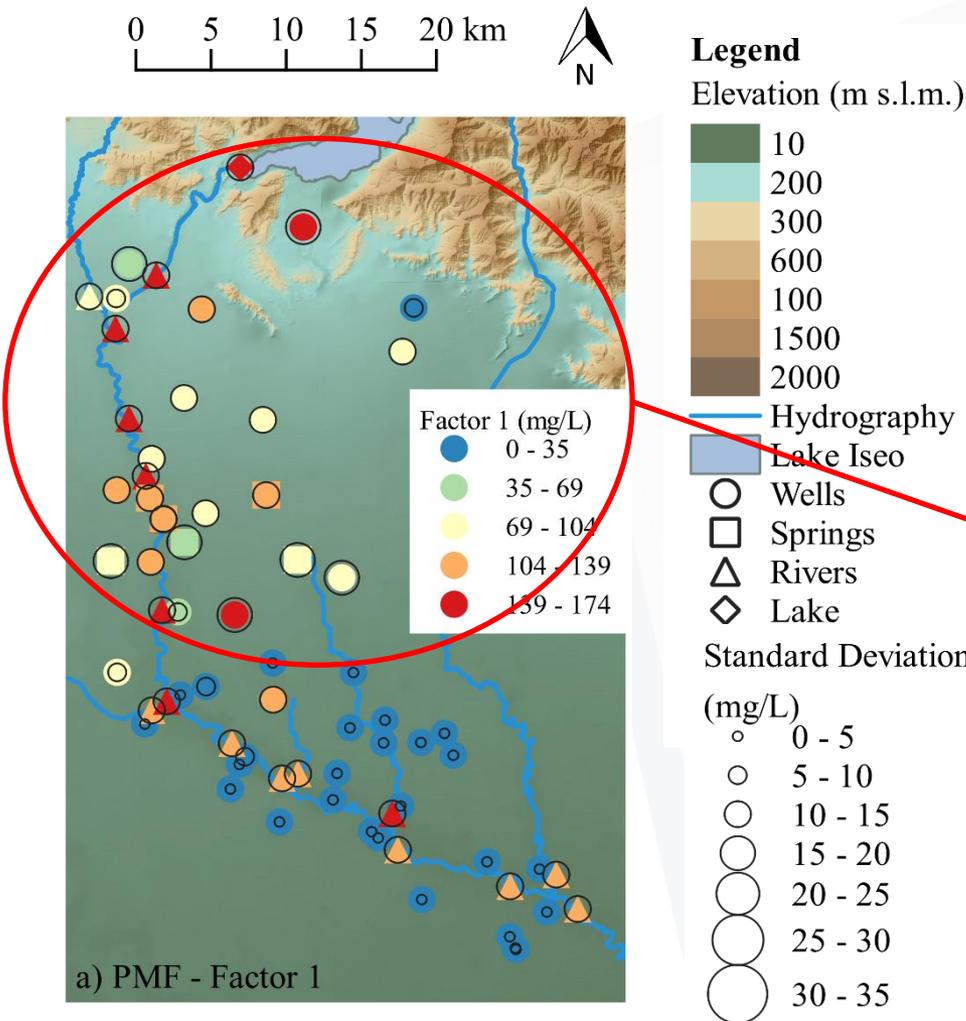
2- **Reducing condition - advanced stages:** As, Fe, P-tot and NH<sub>4</sub>, furthermore it has a significant contribution of Mn

3- **Reducing condition - early stages:** mainly characterized by Mn and SO<sub>4</sub>, with a contribution of the major ions

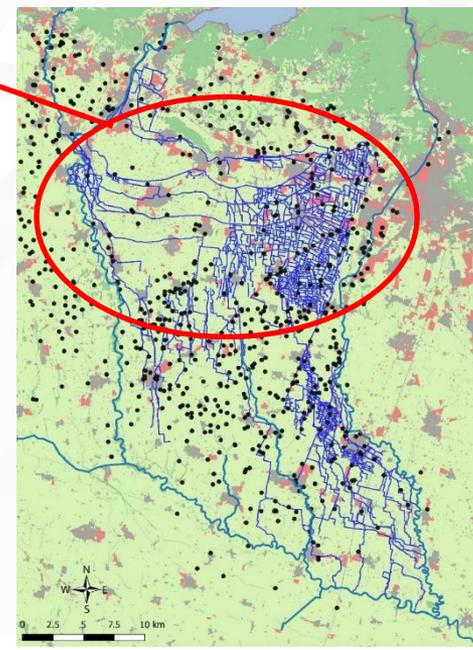
4- **Residence time:** only major ions such as HCO<sub>3</sub>, Cl, Ca, Mg, Na, and K

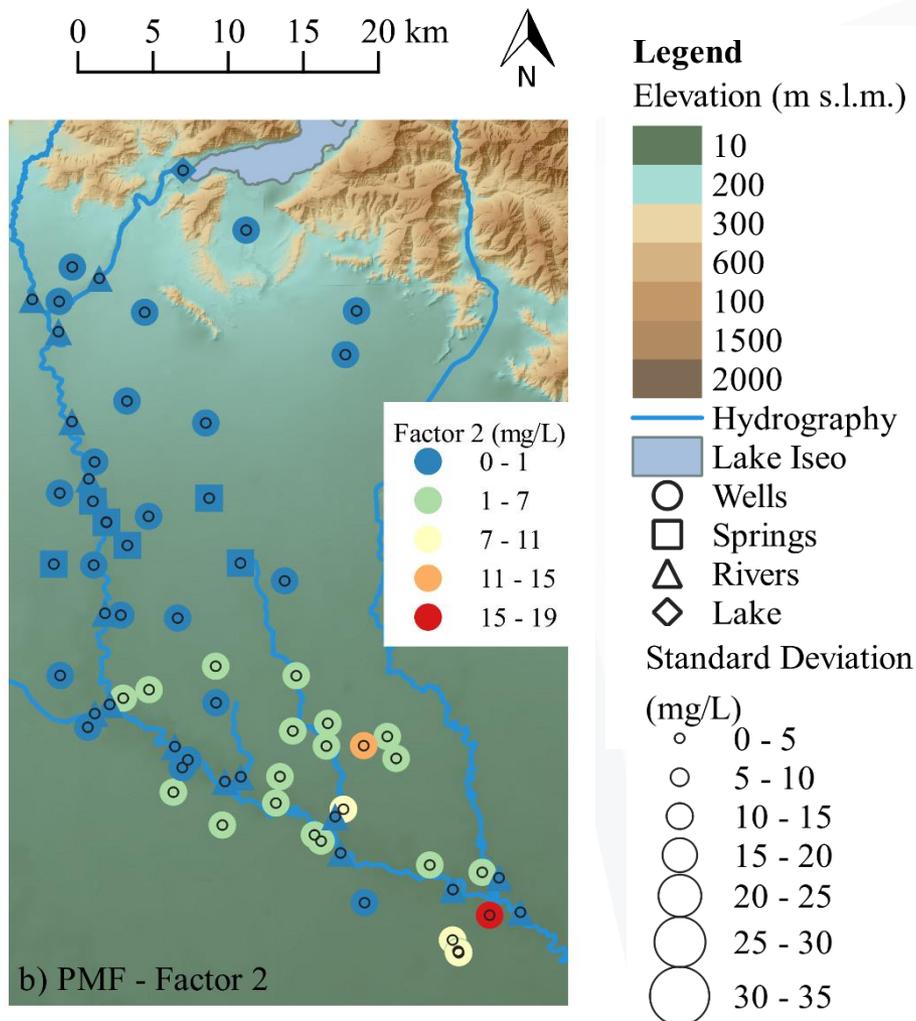
5- **Anthropogenic impact:** NO<sub>3</sub> with contributions also in terms of Cl, SO<sub>4</sub>, Ca and Mg.

# PMF Fac 1 – Lake Iseo and Oglio River



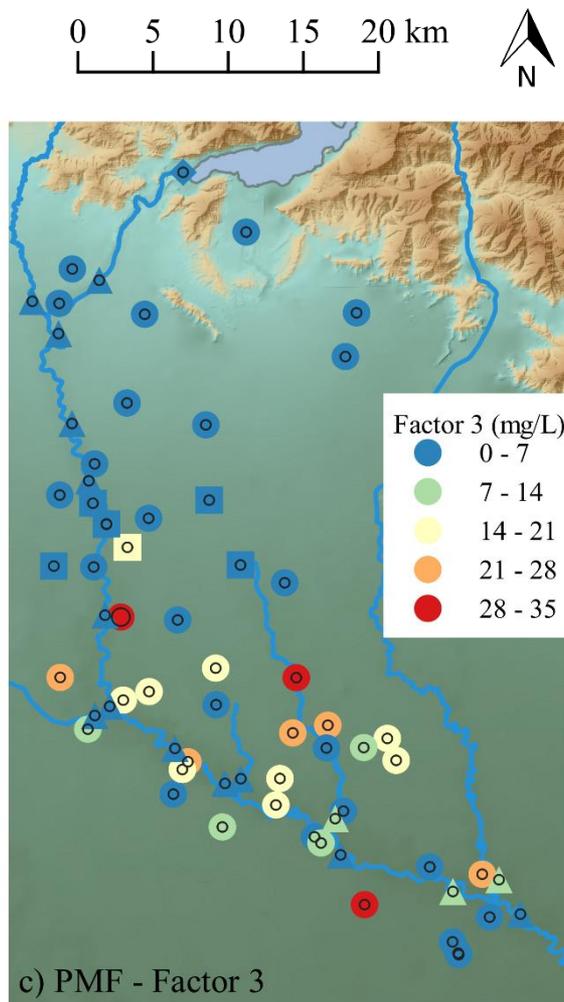
- Water from Lake Iseo flows into the Oglio river
- Higher values in the higher plain wells due to irrigation processes: channels collect water from the Oglio river and spread it over the fields from where it can reach the aquifer





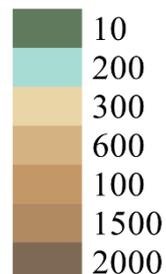
- Sampling points with higher values are wells in the lower plain, where As pollution is stronger
- Groundwater samples in the higher plain, which is unconfined, and surface water samples have no contribution of this factor.

# PMF Fac 3 – Early Red. conditions



## Legend

Elevation (m s.l.m.)



Hydrography

Lake Iseo

Wells

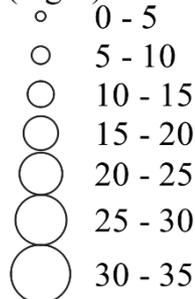
Springs

Rivers

Lake

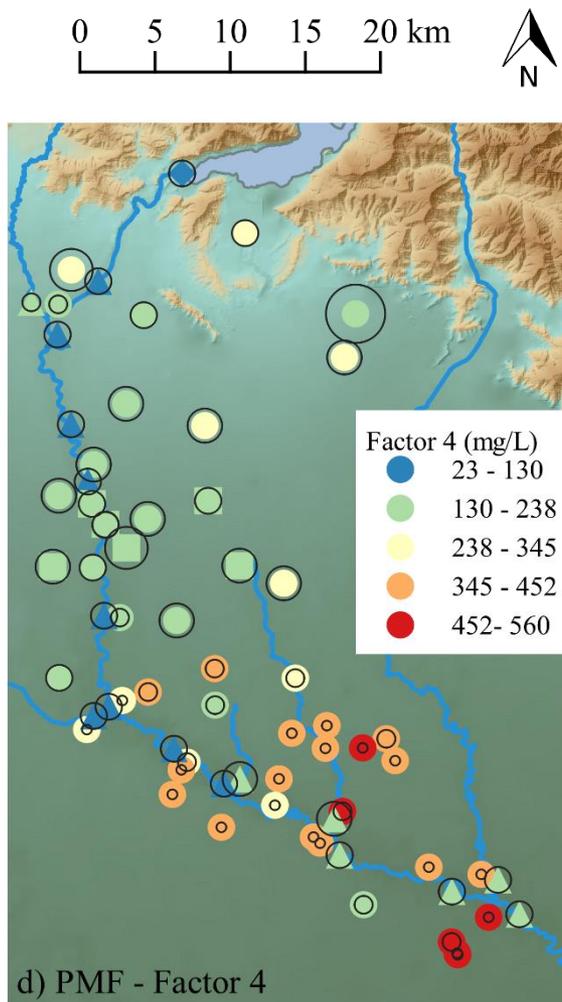
Standard Deviation

(mg/L)



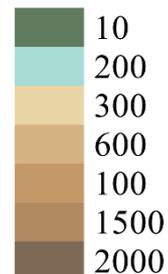
- sampling points with the highest values at the transition between higher and lower plain aquifers
- here groundwater is passing from the oxidising conditions of the higher plain to the reducing conditions of the lower plain, so Mn-oxide reduction can be favoured.
- groundwater samples in the higher plain and surface water samples have no contribution of this factor.

# PMF Fac 4 – Residence time



## Legend

Elevation (m s.l.m.)



Hydrography

Lake Iseo

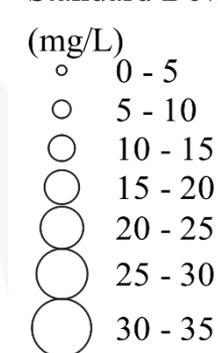
Wells

Springs

Rivers

Lake

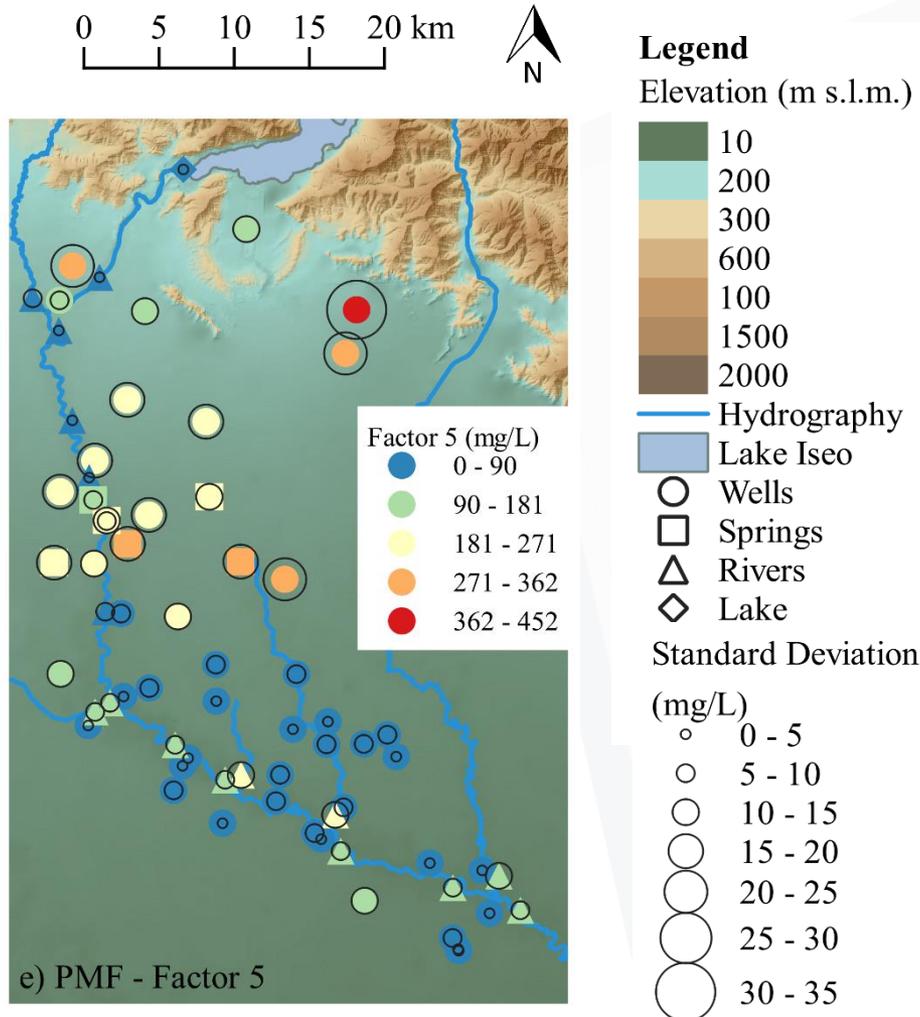
Standard Deviation



➤ Effect of water – rock interactions increasing from north to south, which is the flow direction

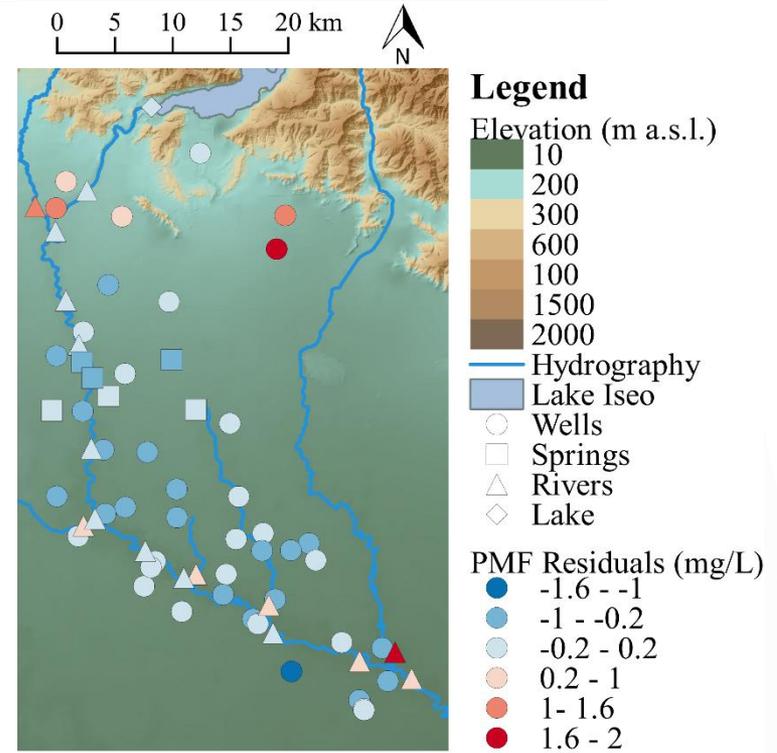
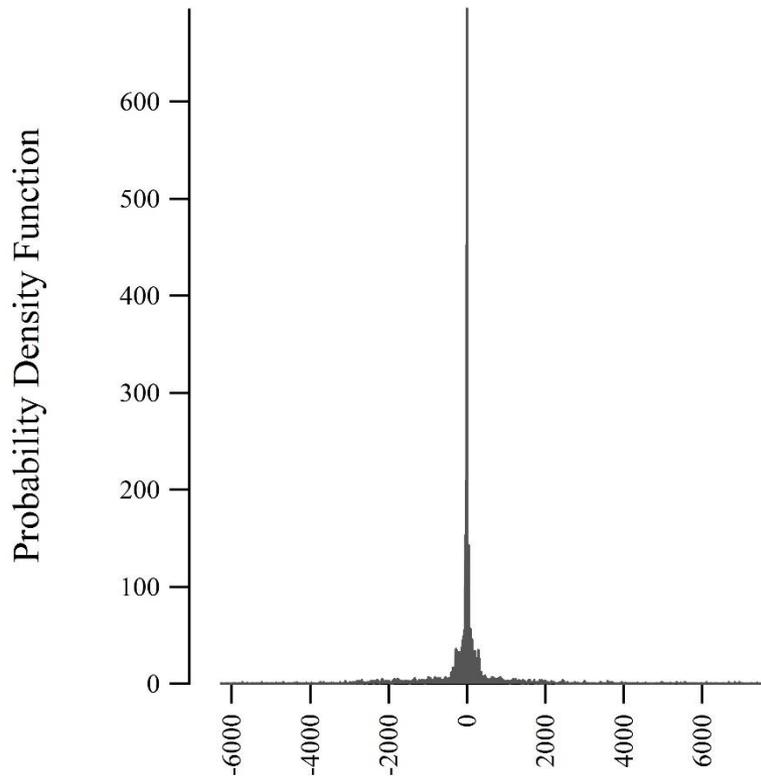
➤ Highlights the gaining behaviour of the river (increasing trend among surface water samples going downstream)

# PMF Fac 5 – Anthropogenic impact



- Anthropogenic impact is higher in the higher plain (unconfined aquifer is more vulnerable and there is no reduction of the nitrate)
- This factor highlights the gaining behaviour of the river and the impact of tributaries

# PMF Residual analysis

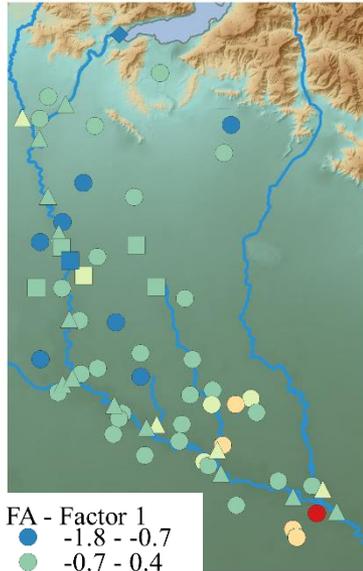


Residuals have a symmetric distribution

Residuals do not have a spatial pattern

# COMPARISON WITH FA

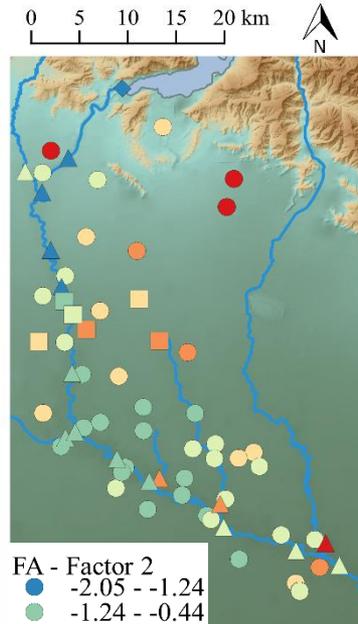
Factor analysis output consist in 3 factors:



- FA - Factor 1
- -1.8 - -0.7
  - -0.7 - 0.4
  - 0.4 - 1.5
  - 1.5 - 2.6
  - 2.6 - 3.7
  - 3.7 - 4.9

1- Reducing condition advanced + early stages

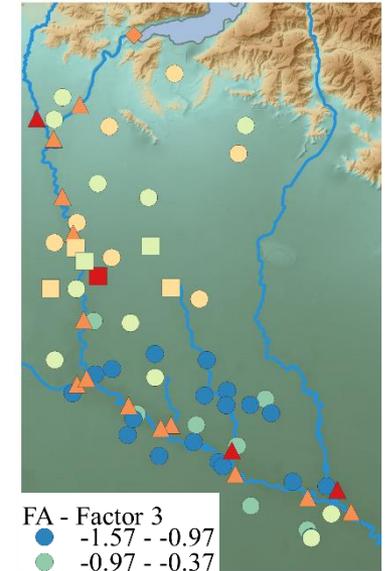
Ptot, NH<sub>4</sub>, Na, As, Fe, Mn  
SO<sub>4</sub>



- FA - Factor 2
- -2.05 - -1.24
  - -1.24 - -0.44
  - -0.44 - 0.37
  - 0.37 - 1.17
  - 1.17 - 1.98
  - 1.98 - 2.78

2- Anthropogenic impact + water rock interactions

HCO<sub>3</sub>, Cl, NO<sub>3</sub>, Ca,  
Mg



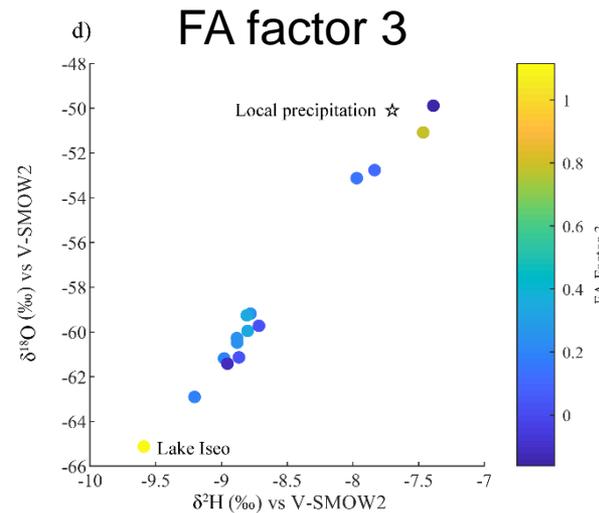
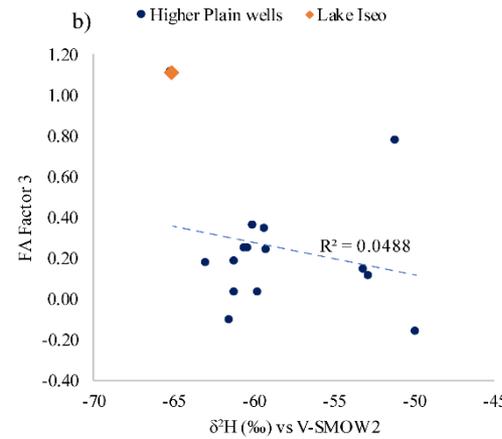
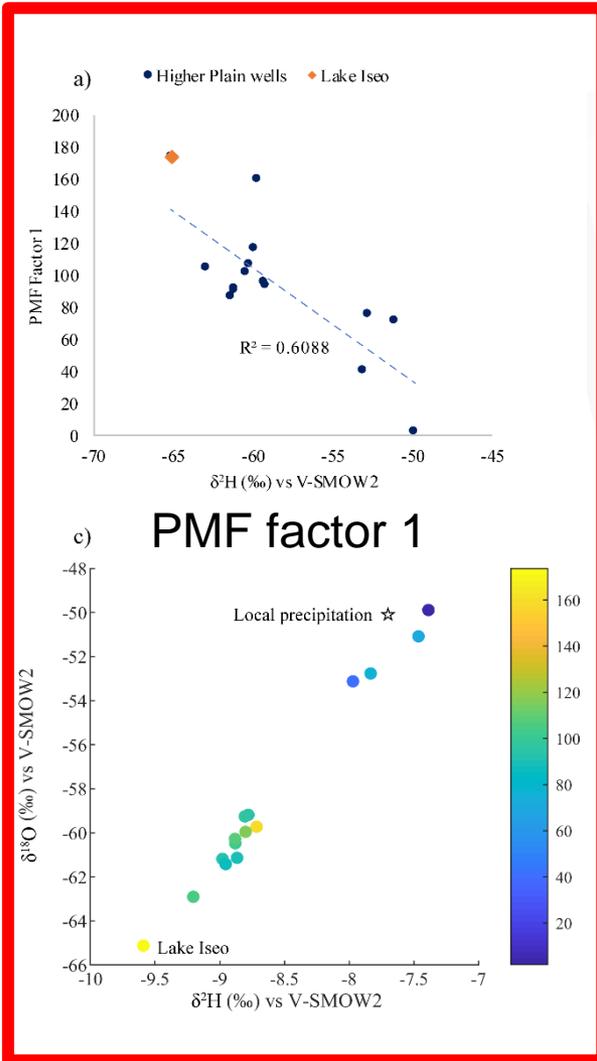
- FA - Factor 3
- -1.57 - -0.97
  - -0.97 - -0.37
  - -0.37 - 0.22
  - 0.22 - 0.82
  - 0.82 - 1.42
  - 1.42 - 2.02

3- Surface water

O<sub>2</sub>, Cl, NO<sub>3</sub>, SO<sub>4</sub>, NH<sub>4</sub>  
K

# COMPARISON WITH FA - validation

Validation of PMF factor 1 and FA factor 3 (surface water) with isotopic data



Two main sources of recharge in the higher plain aquifer:

- **Surface water from irrigation:** more negative values ( $\leftarrow$  precipitation from higher latitude and altitude)
- **Local precipitation:** less negative values ( $\leftarrow$  precipitation from local latitude and altitude)

Higher plain well samples more recharged by irrigation have isotopic signature closer to the signature of Lake Iseo

**PMF fac 1 shows higher correlation with isotopic data**



- Multivariate analysis allows for a synthetic description of the system
- All the main phenomena characterizing the system are identified and quantified
- PMF leads to a more environmentally interpretable solution as compared to FA



- PMF could be a suitable tool to perform exploratory analysis, supporting the development of the conceptual model



Water Research 159 (2019) 122–134



Contents lists available at [ScienceDirect](#)

Water Research

journal homepage: [www.elsevier.com/locate/watres](http://www.elsevier.com/locate/watres)



Groundwater and surface water quality characterization through positive matrix factorization combined with GIS approach



C. Zanotti <sup>a,\*</sup>, M. Rotiroti <sup>a</sup>, L. Fumagalli <sup>a</sup>, G.A. Stefania <sup>a</sup>, F. Canonaco <sup>b</sup>, G. Stefenelli <sup>b</sup>, A.S.H. Prévôt <sup>b</sup>, B. Leoni <sup>a</sup>, T. Bonomi <sup>a</sup>

<sup>a</sup> Department of Earth and Environmental Sciences, University of Milano-Bicocca, Piazza della Scienza, 1, 20126, Milano, Italy

<sup>b</sup> Laboratory of Atmospheric Chemistry, Paul Scherrer Institute, 5232, Villigen-PSI, Switzerland

**For more information**

Author contact: [chiara.zanotti@unimib.it](mailto:chiara.zanotti@unimib.it)

full paper at: <https://doi.org/10.1016/j.watres.2019.04.058>