# Presentation Outline

- **Introduction**

- **Study area description**

- **Data Set**

  - **General description, descriptive statistics, histogram, boxplot**

- **Methodology used**

  - **Scenarios formulation**

- **Results & Discussions**

- **Conclusion**

# Introduction

- **Urban groundwater (GW) is usually vulnerable to pollution.**

- **The main sources of GW quality degradation are:**

  - **Anthropogenic activities**
  - **Natural processes**
  - **Atmospheric input**

- **Subdividing the region into zones based on GW quality is usually undertaken.**

- **Recently, several methods are used, such as Genetic Algorithm, Model-Based Approach, Bayesian Approach, cluster analysis … etc.**

- **In this study, spatially constrained multivariate clustering method (SCMC) is used to subdivide Madina city (West KSA) into several zones based on six GW chemicals.**
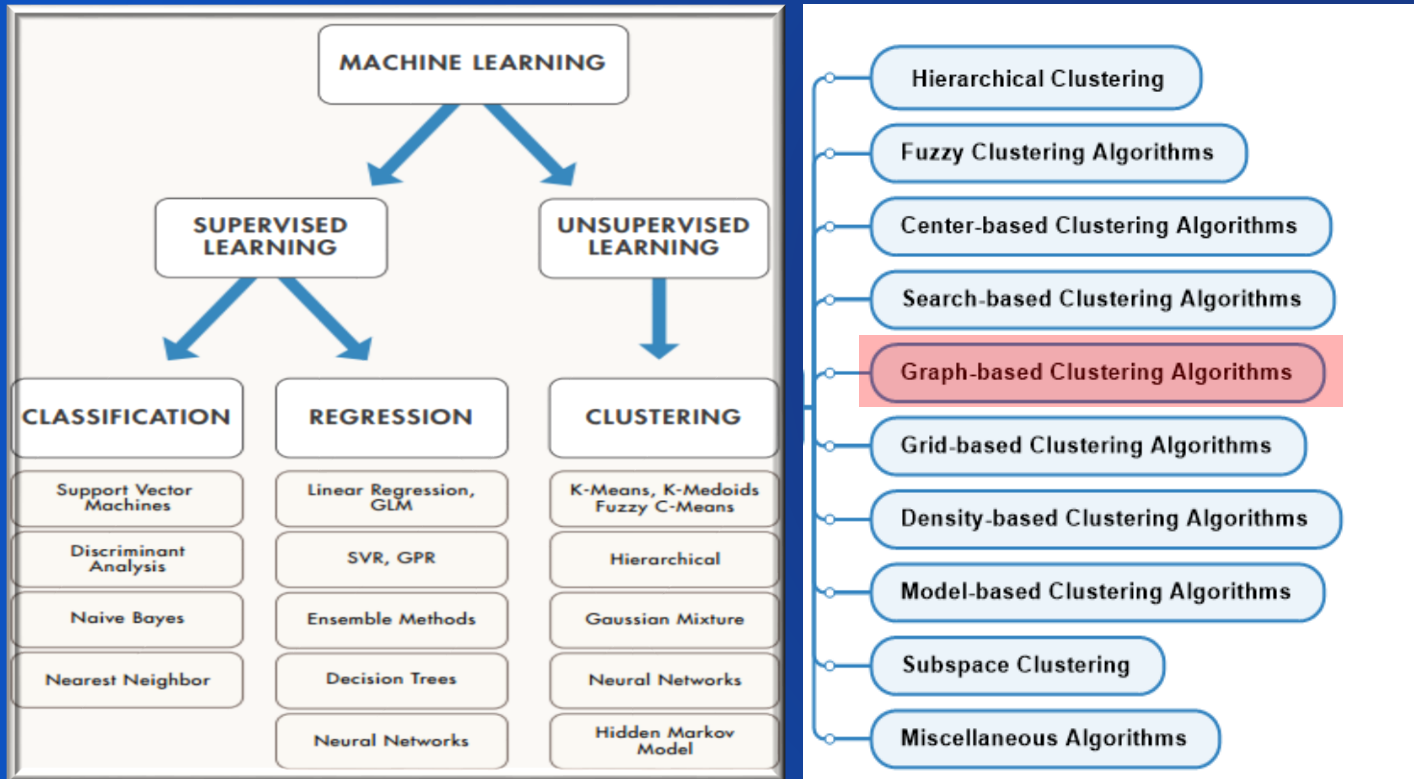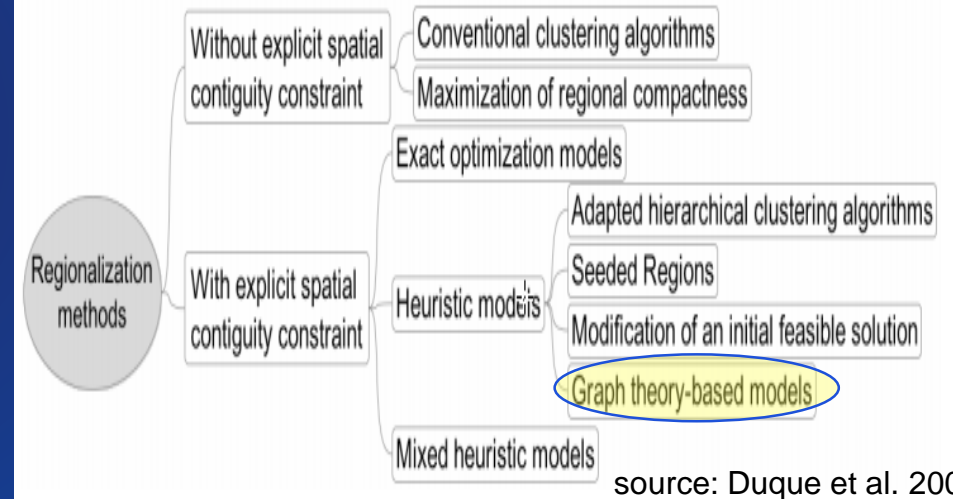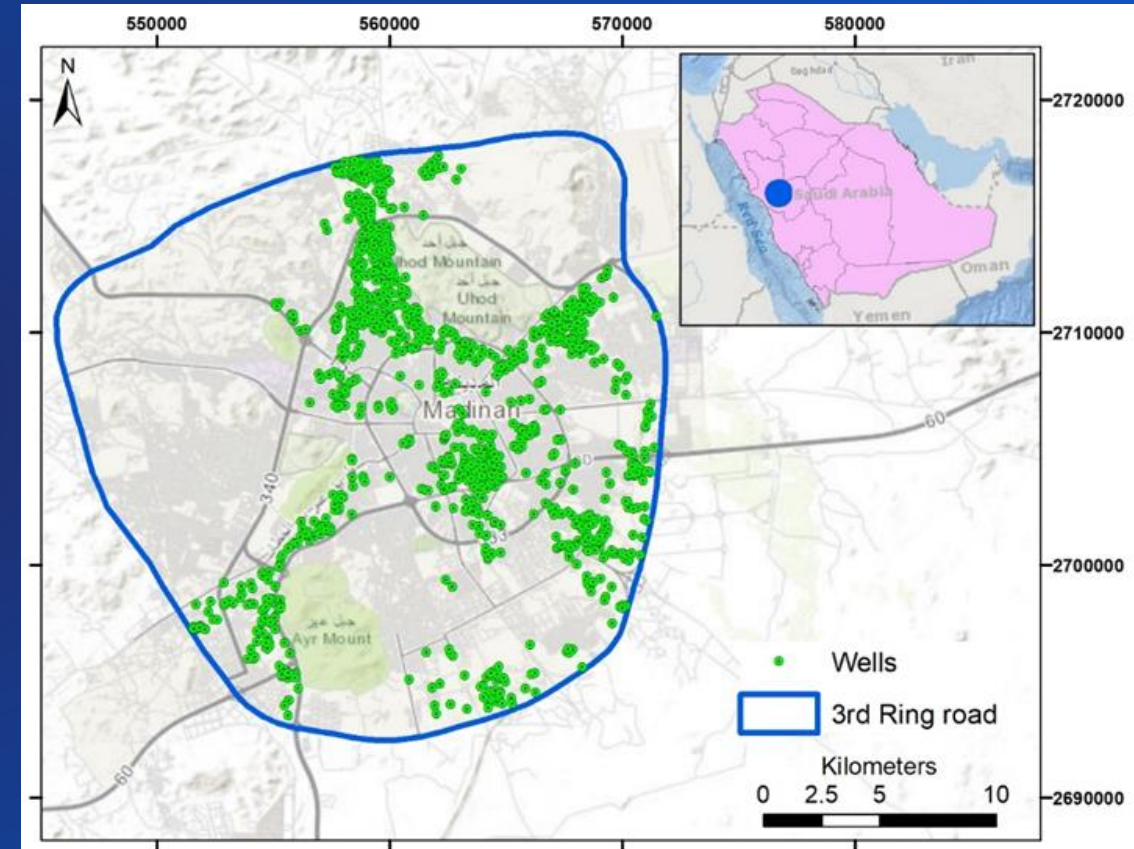
# Introduction (cluster analysis overview)



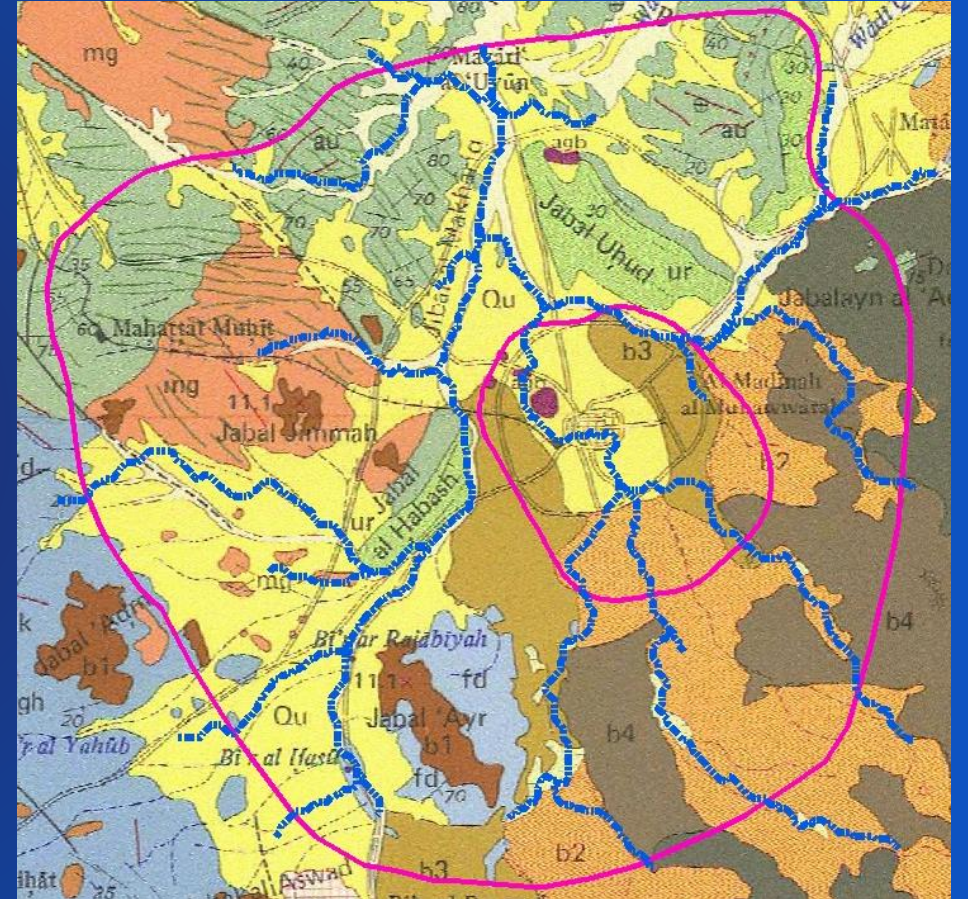Figure 1. Taxonomy of methods for solving regionalization problems

source: Duque et al. 2007

# Study area description (General Location)

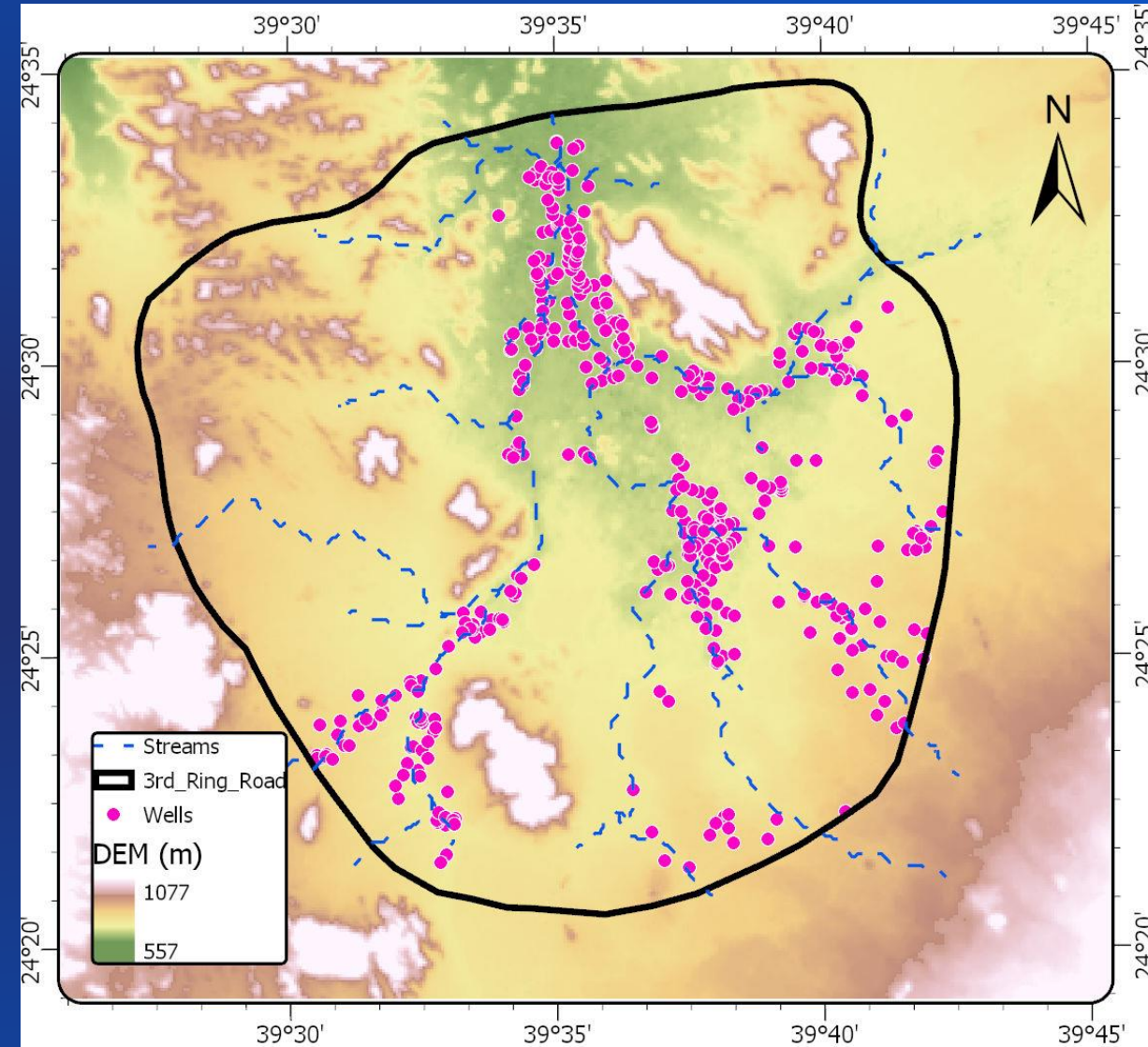❖ **Madinah city (West KSA) is selected.**

❖ **The population about 1.25 million + 10 million annual visitors.**

❖ **covering an area of 522 km$^2$** .

❖ **In this study, The wells inside the 3rd ring road is selected for analysis.**

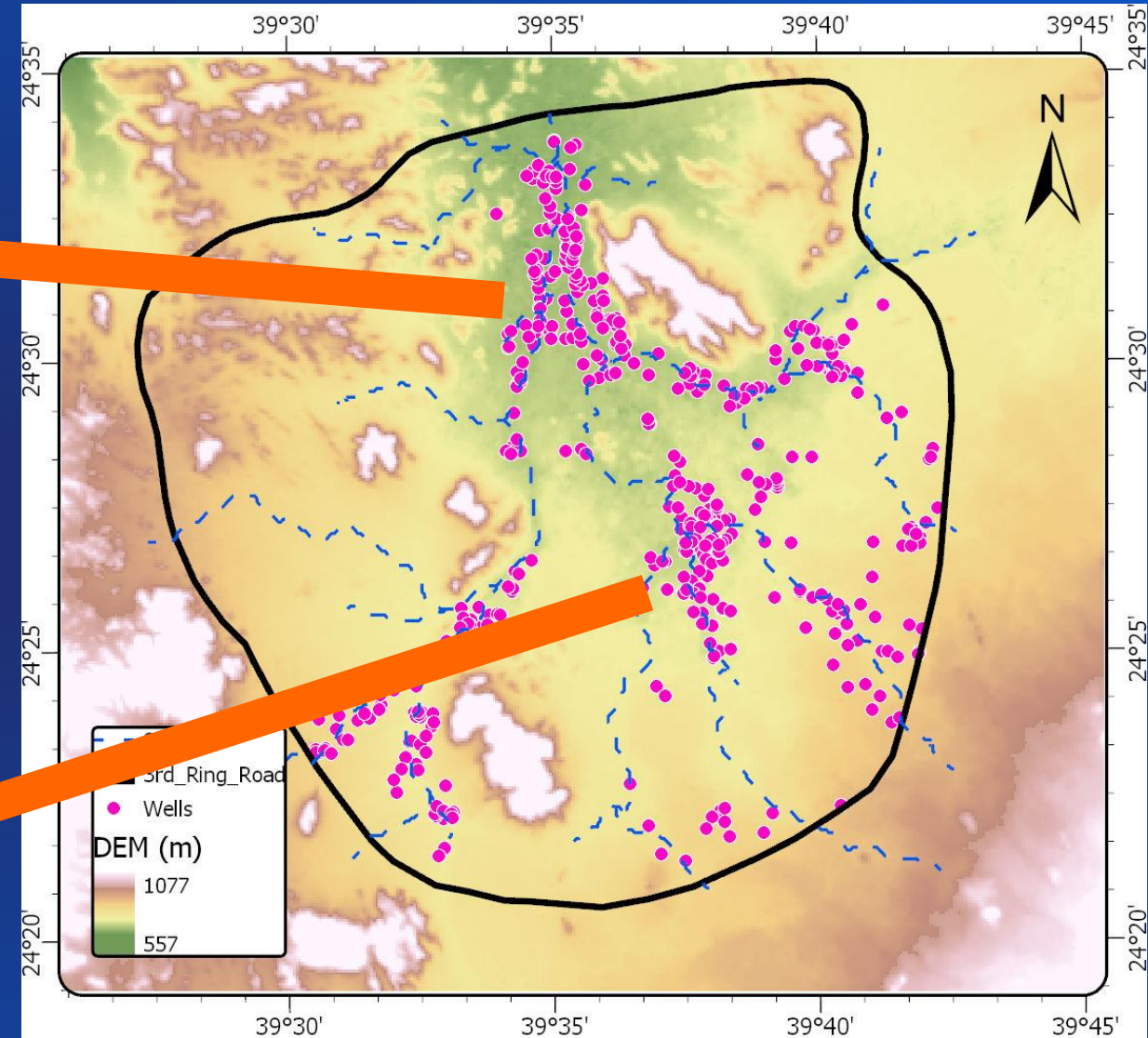❖ **These wells are located in private farms.**

# Study area description (Geology)

❖ **The geology of the study area consists mainly of three parts;**

  ❖ **lava plateaus (volcanic basalt flows)**

  ❖ **alluvial deposits**

  ❖ **rock outcrops (pre-cambrian rock).**

❖ **The first two parts are the places of shallow groundwater aquifers.**

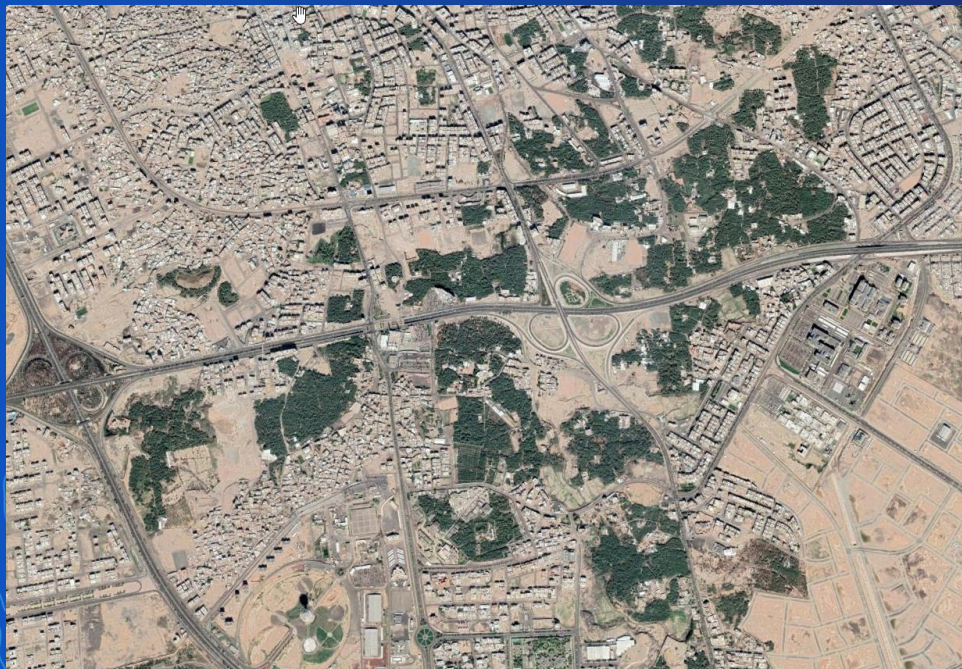❖ **50% of the study area is covered by volcanic basalt rocks**

# Study area description (Topography)

- **The elevation ranges from 570 m (a.m.s.l.) up to 1,100 m.**

- **Strong relationship between wells location and watercourses (ephemeral streams)**
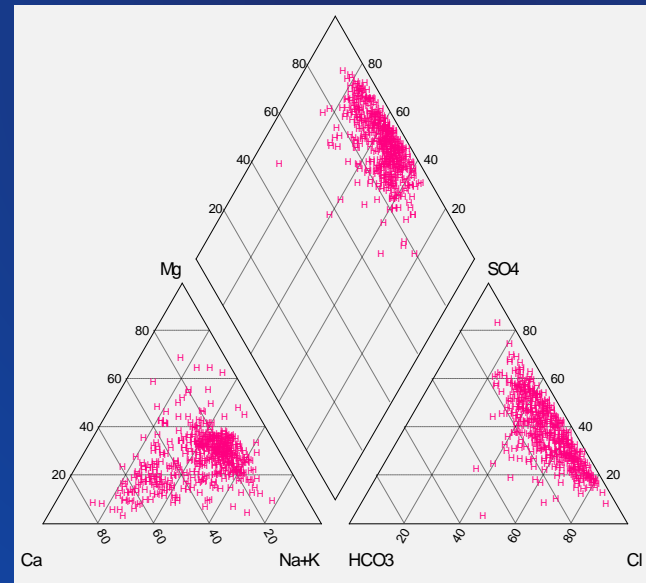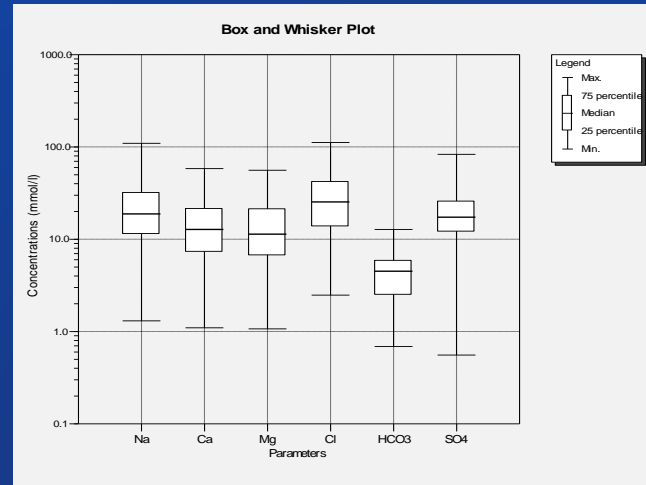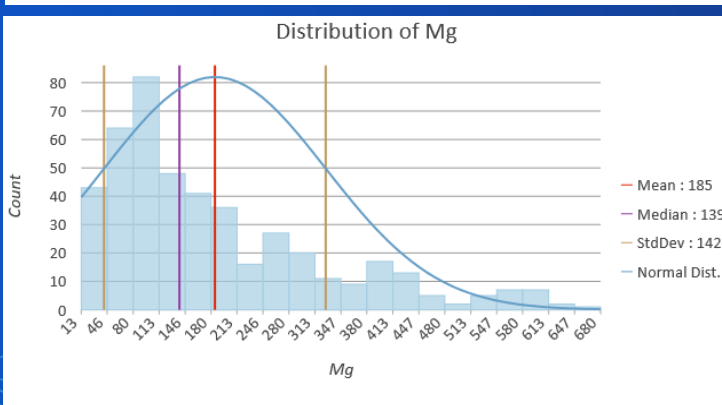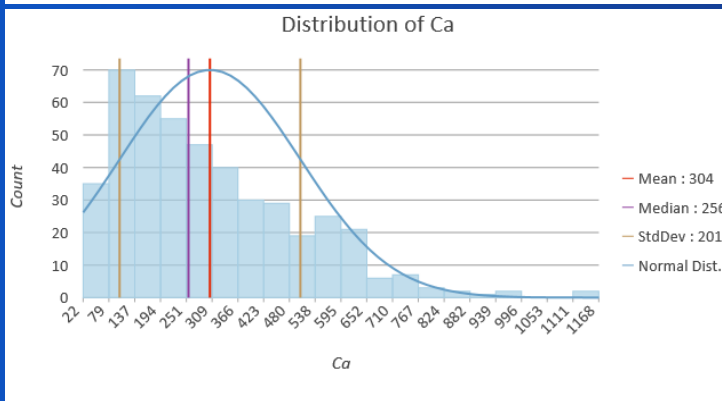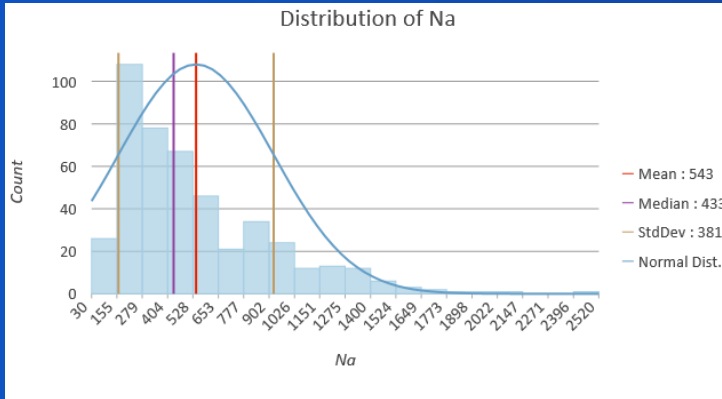
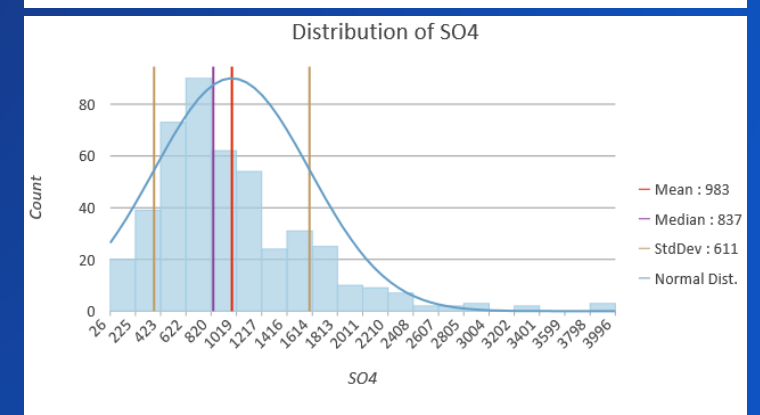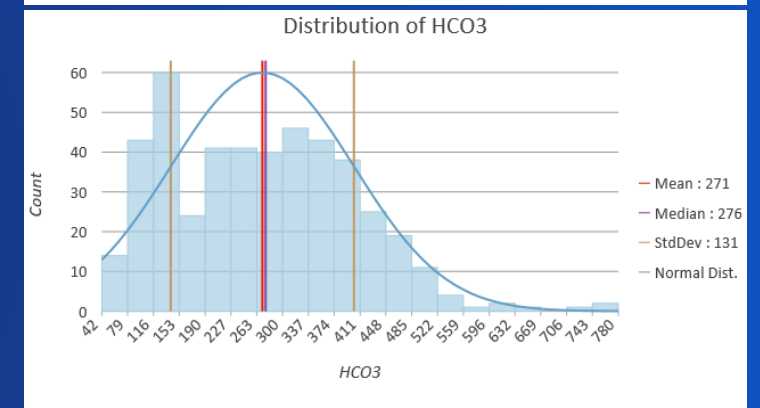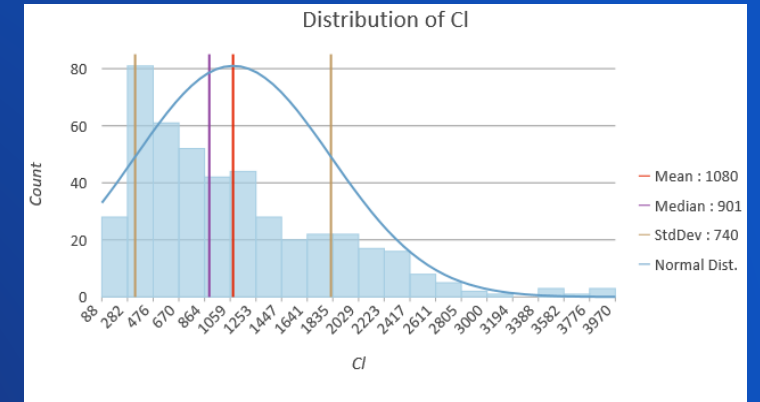# Study area description (Wells location)

# Data set

➢ **456 private farms inside the 3$^{rd}$ ring road are visited.**

➢ **From each farm, one well location is registered using GPS.**

➢ **Water samples from the wells are collected and taken to the laboratory for analysis (pH, TDS, EC, hardness, turbidity, alkalinity, color, ions (_cations and anions_).**

➢ **Three cations (Na, Ca, Mg) and three anions (Cl, HCO$_3$, SO$_4$) are selected for cluster analysis.**

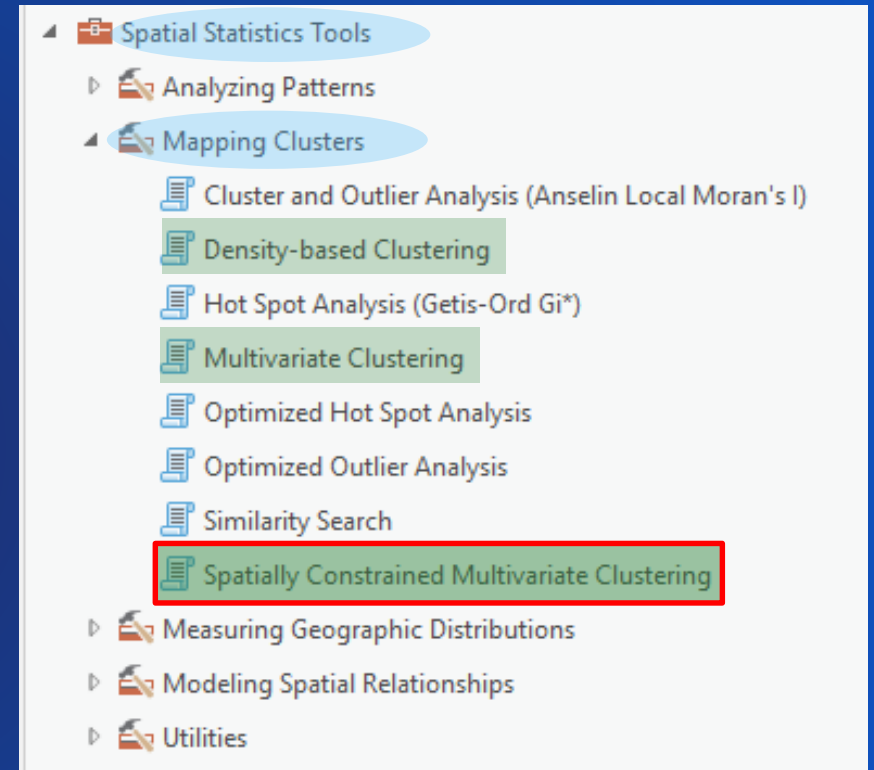| | Na | Ca | Mg | Cl | HCO$_3$ | SO$_4$ |
|---|---|---|---|---|---|---|
| Min. | 30 | 22 | 13 | 88 | 42 | 26 |
| Max. | 2520 | 1168 | 680 | 3970 | 780 | 3996 |
| Avg. | 543 | 304 | 185 | 1080 | 271 | 982 |
| Media | 433 | 256 | 139 | 901 | 276 | 836 |
| STD | 381 | 201 | 142 | 740 | 131 | 611 |

# Data set



Piper plots (Trilinear diagram)

# Methodology (clustering in ArcGIS Pro)
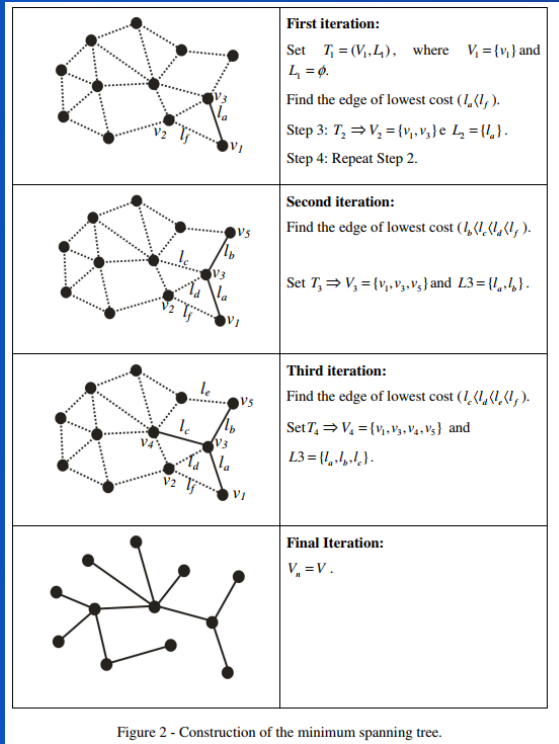
- **The Three new cluster methods in ArcGIS Pro are:**
  - **density-based clustering,**
  - **multivariate clustering and**
  - **spatially constrained multivariate clustering (SCMC).**

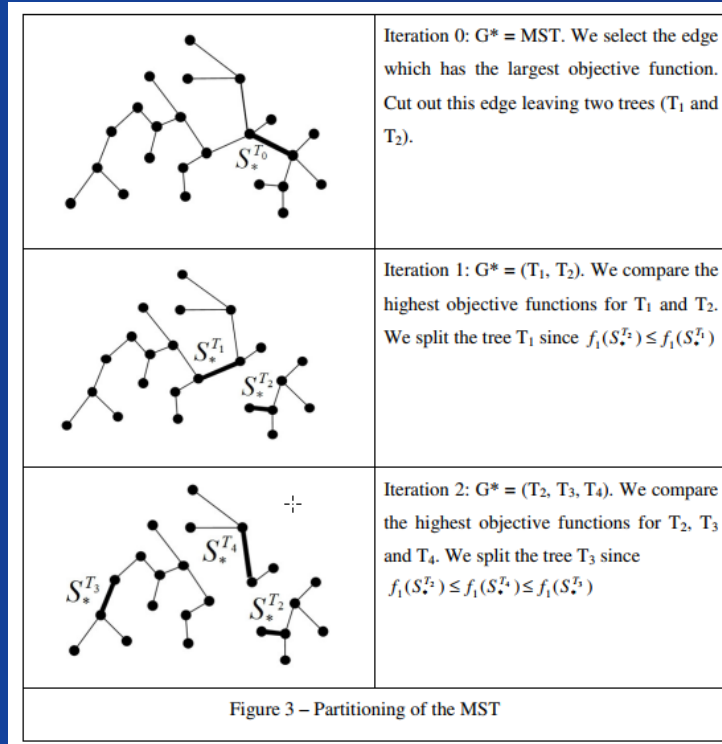SCMS is the process of grouping of the observations based on the attributes similarity and location similarity using multiple objective optimization. maximizing within-group similarity
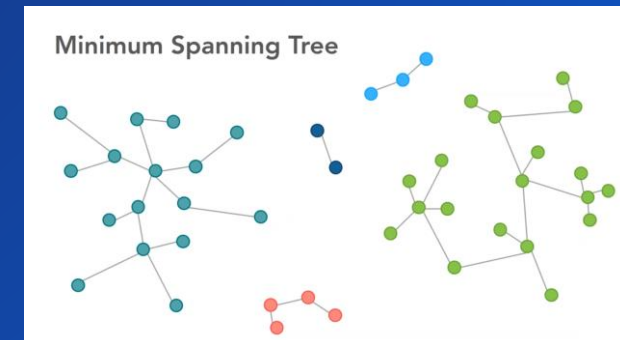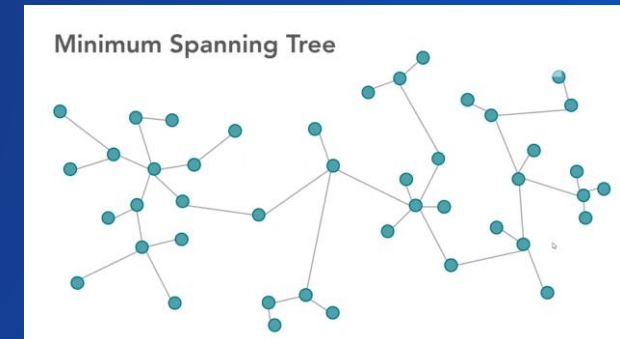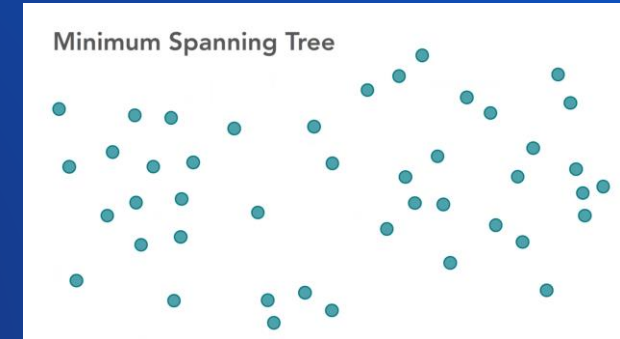
# Methodology (theoretical background)

- **In ArcGIS Pro, spatially constrained multivariate clustering tool uses Spatial K'luster Analysis by Tree Edge Removal (SKATER) algorithm which is based on minimum spanning tree (MST) method**



Figure 2 - Construction of the minimum spanning tree.



Figure 3 – Partitioning of the MST



construct the network graph by connecting the contiguous nodes with lowest cost.

find the shortest path that minimizes the sum of dissimilarity (or maximizing the sum of similarities) – minimum spanning tree

Source: Assuncao et al (2006)

# Methodology

- **Spatially constrained Multivariate Clustering (SCMC) methods in ArcGIS Pro has two main groups of input parameters:**

- **Three required input:**
  - Input layer
  - The name of output layer
  - Selected attributes for analysis

- **Five optional input:**
  - **Cluster size constraints** (None, No. of features, Attribute value)
  - Number of clusters
  - **Spatial constraints**
  - Permutations Membership probabilities
  - Output table for evaluating number of clusters

# Methodology (Scenarios formulation)

- **Four groups of scenarios are developed based on the optional input parameters:**
  - **Group (A) scenario** : no optional input and the optimum number of clusters is computed automatically.
  - **Group (B) scenarios** : three optimum No. of clusters are specified (2, 3 and 7).
  - **Group (C) scenarios**: min. No. of features per cluster is specified (20, 40), with fixed number of cluster (= 7 clusters).
  - **Group (D) scenarios**: min. no. and max. no. of features per clusters are specified ("25, 150" and "50, 100"), the optimum number of clusters is computed automatically.

**Group (A) scenario**

Cluster Size Constraints
None
Number of Clusters
Spatial Constraints
Trimmed Delaunay triangulation
Permutations to Calculate Membership Probabilities
Output Table for Evaluating Number of Clusters

**Group (B) scenarios**

Cluster Size Constraints
None
Number of Clusters          3
Spatial Constraints
Trimmed Delaunay triangulation
Permutations to Calculate Membership Probabilities          1000
Output Table for Evaluating Number of Clusters

**Group (C) scenarios**

Cluster Size Constraints
Number of features
Minimum per Cluster          40
Number of Clusters          7
Spatial Constraints
Trimmed Delaunay triangulation
Permutations to Calculate Membership Probabilities          1000
Output Table for Evaluating Number of Clusters

**Group (D) scenarios**

Cluster Size Constraints
Number of features
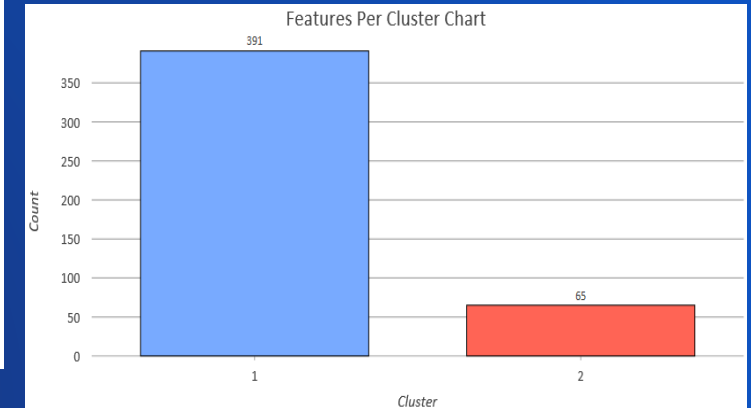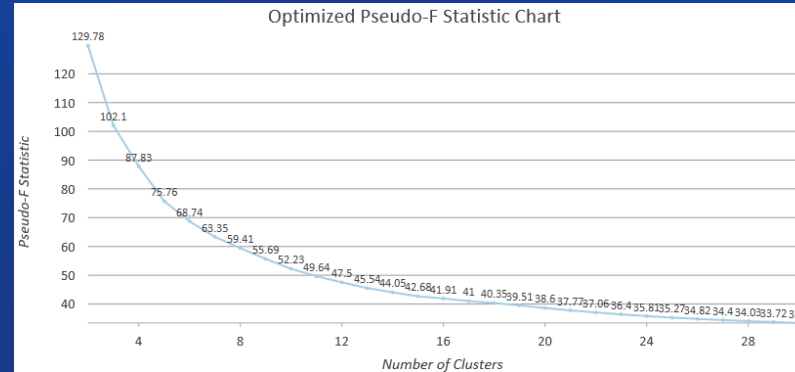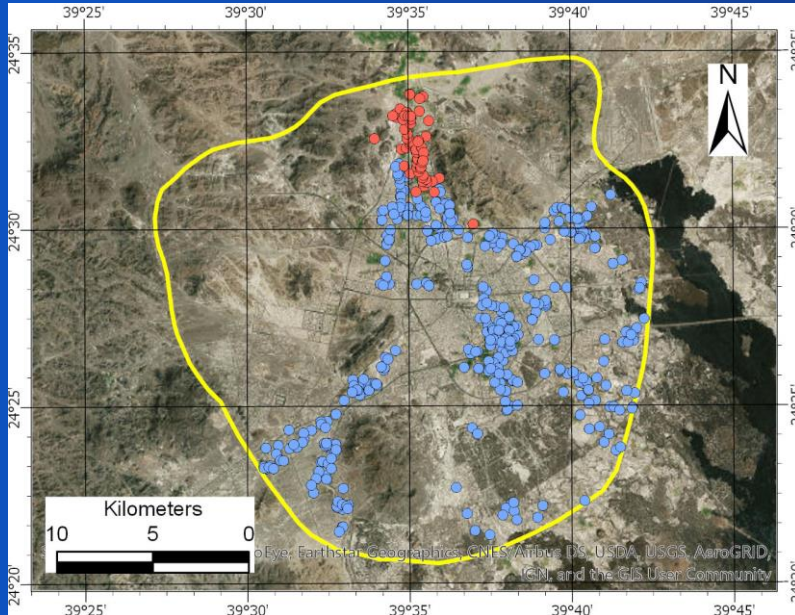Minimum per Cluster          25
Fill to Limit          150
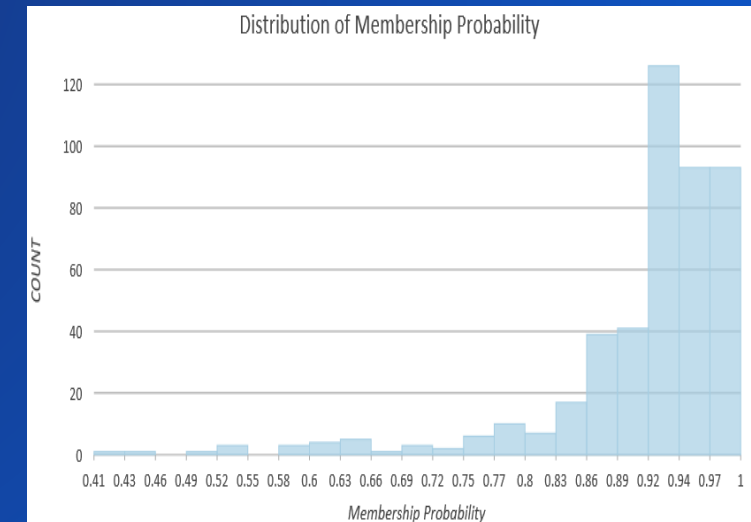Spatial Constraints
Trimmed Delaunay triangulation
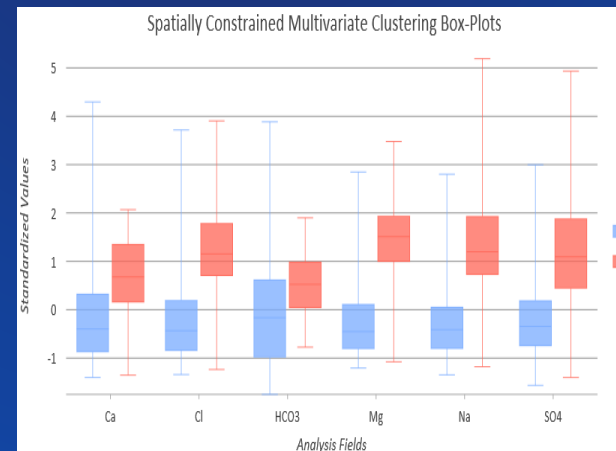
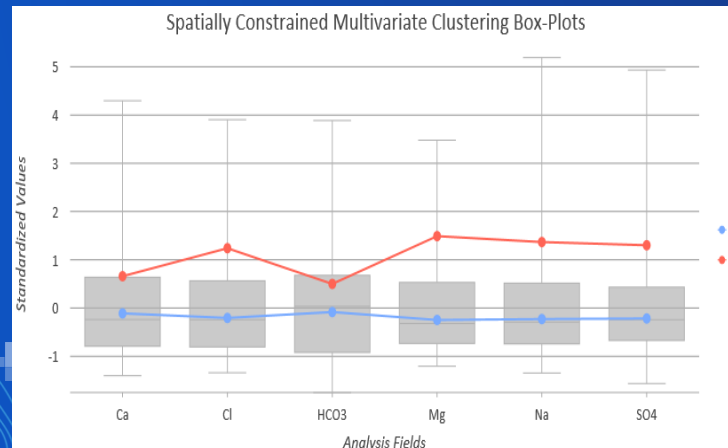# Results & Discussions (Group (A) scenario)

No optional input and the optimum number of clusters is computed automatically
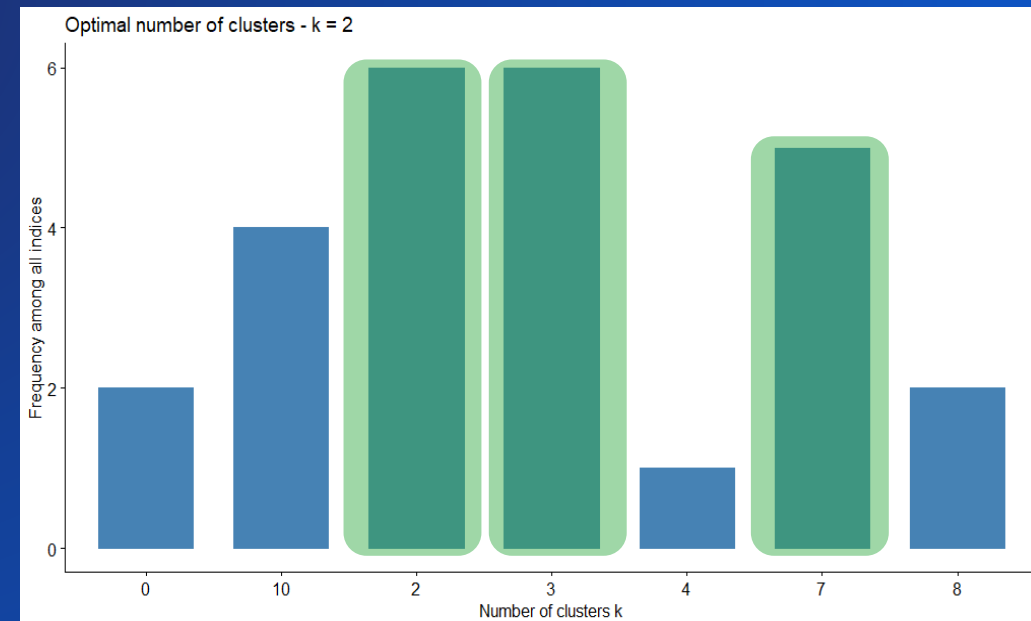
# Results & Discussions (Group (B) scenarios)

- **Most of the clusters methods need from the user to specify the optimum number of clusters**
- **Unfortunately, this is still unsolved problem and there is no definitive answer to this question.**
- **Determining the optimal number of clusters is somehow subjective.**
- **In this study, the optimum number of clusters is determined by evaluating 30 indices using R programming language (NbClust R package).**
- **Selection of the optimum No. of cluster is based on the "majority rule".**

```
Among all indices:
===================
* 2 proposed   0 as the best number of clusters
* 6 proposed   2 as the best number of clusters
* 6 proposed   3 as the best number of clusters
* 1 proposed   4 as the best number of clusters
* 5 proposed   7 as the best number of clusters
* 2 proposed   8 as the best number of clusters
* 4 proposed  10 as the best number of clusters
```
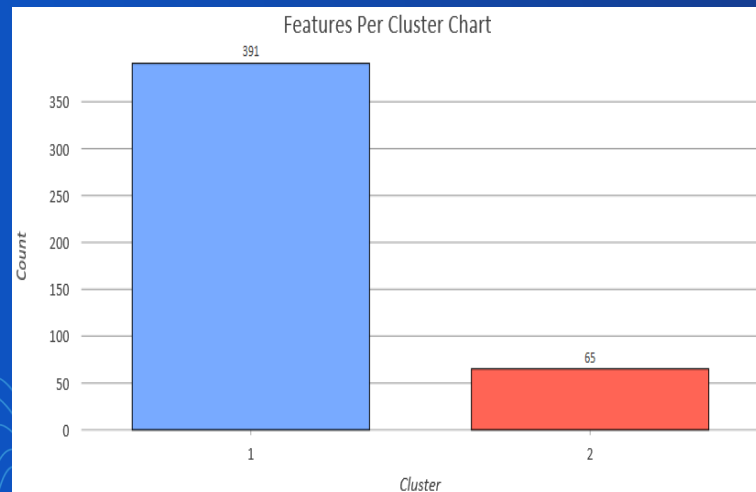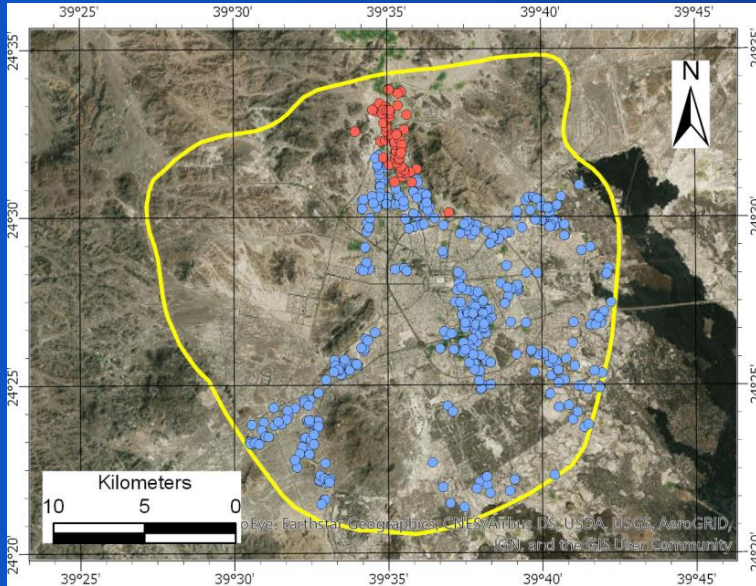
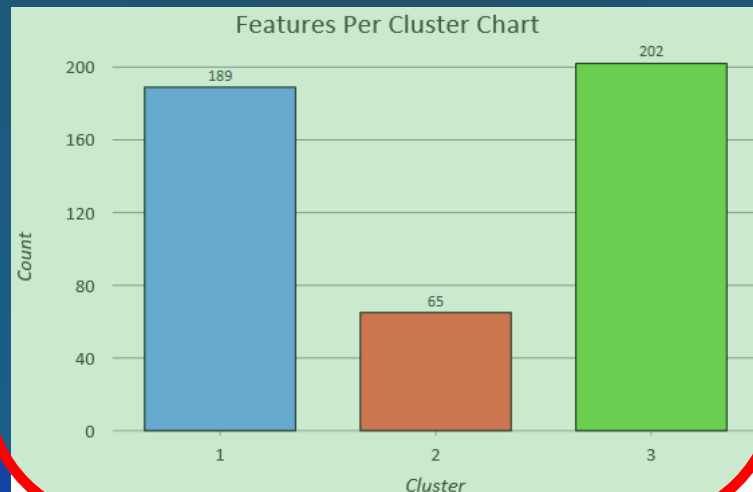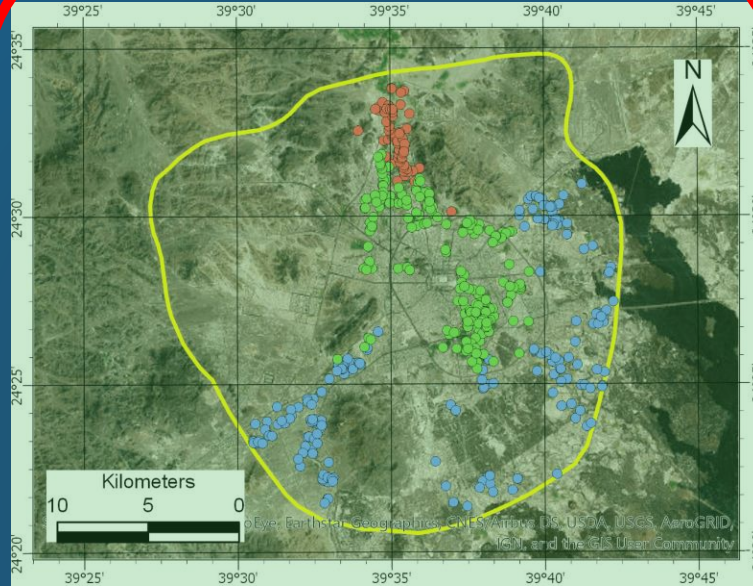- **2, 3 and 7 are selected as the optimum number of clusters**


Optimal number of clusters - k = 2

# Results & Discussions (Group (B) scenarios)

## Three optimum No. of clusters are specified (2, 3 and 7)



No. of clusters = 2

No. of clusters = 3

No. of clusters = 7

# Results & Discussions (Group (B) scenarios)

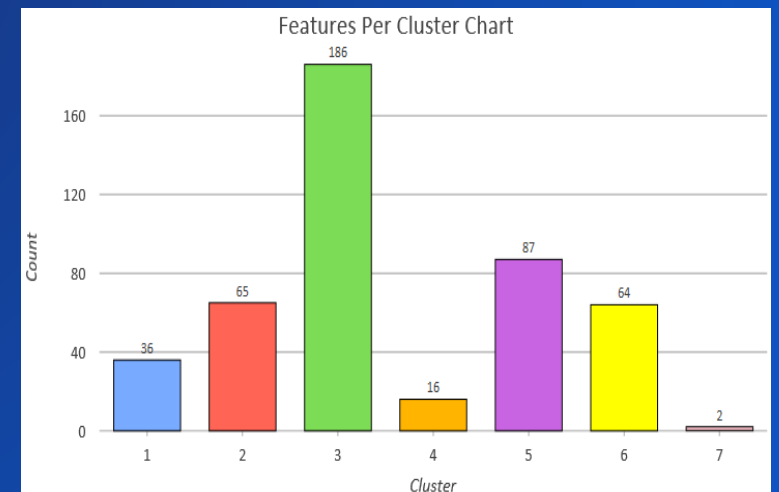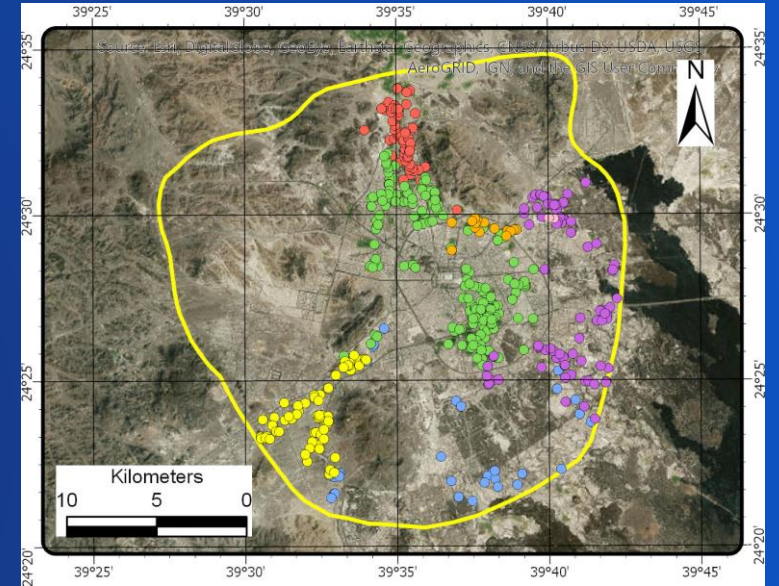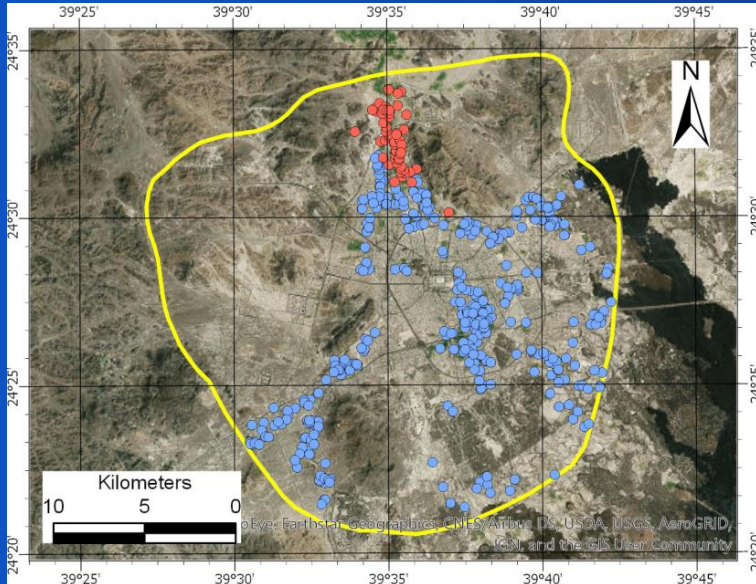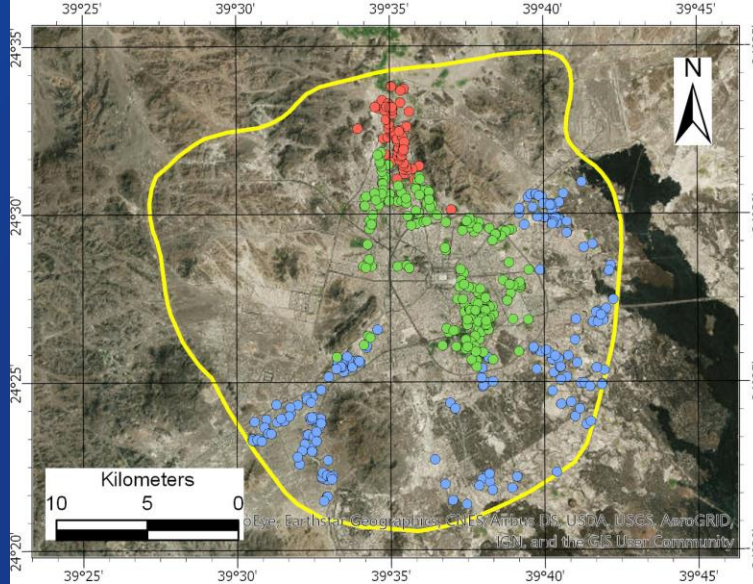# Results & Discussions (Group (B) scenarios)



No. of clusters = 2

No. of clusters = 3

No. of clusters = 7

# Results & Discussions (Group (B) scenarios)



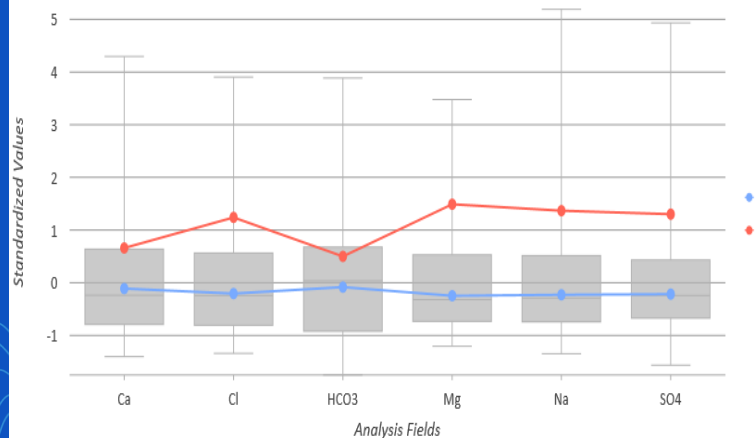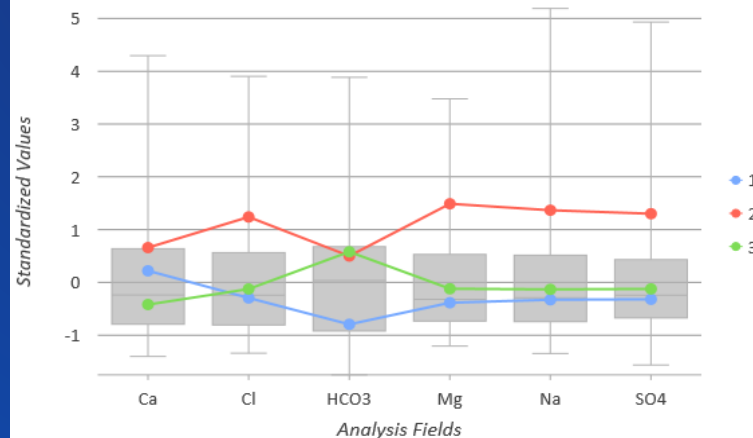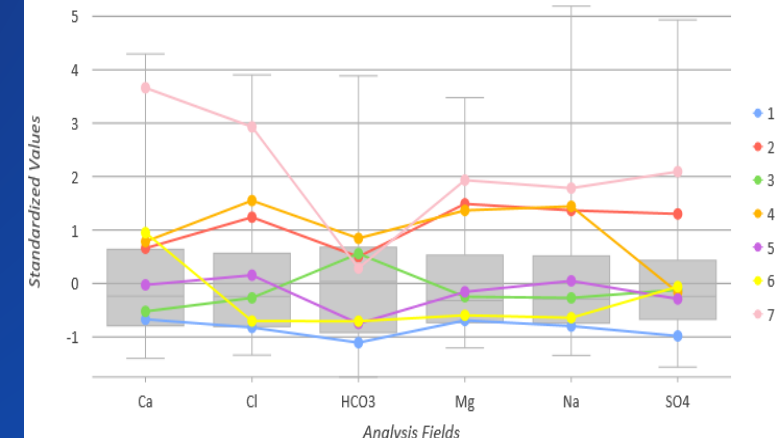| No. of clusters = 2 | No. of clusters = 3 | No. of clusters = 7 |

```
Variable      Mean    Std. Dev.       Min         Max        R2
     MG  184.589912 142.241319 13.000000   680.000000 0.369584
     NA  542.953947 380.473044 30.000000  2520.000000 0.311333
    SO4  982.528509 610.161313 26.000000  3996.000000 0.282234
     CL 1079.949561 739.258629 88.000000  3970.000000 0.256788
     CA  303.885965 200.843004 22.000000  1168.000000 0.072334
   HCO3  271.379386 130.798051 42.000000   780.000000 0.041548
```

```
Variable      Mean    Std. Dev.       Min         Max        R2
   HCO3  271.379386 130.798051 42.000000   780.000000 0.446638
     MG  184.589912 142.241319 13.000000   680.000000 0.384542
     NA  542.953947 380.473044 30.000000  2520.000000 0.319470
    SO4  982.528509 610.161313 26.000000  3996.000000 0.290497
     CL 1079.949561 739.258629 88.000000  3970.000000 0.262823
     CA  303.885965 200.843004 22.000000  1168.000000 0.160296
Writing Results to Output Table
```
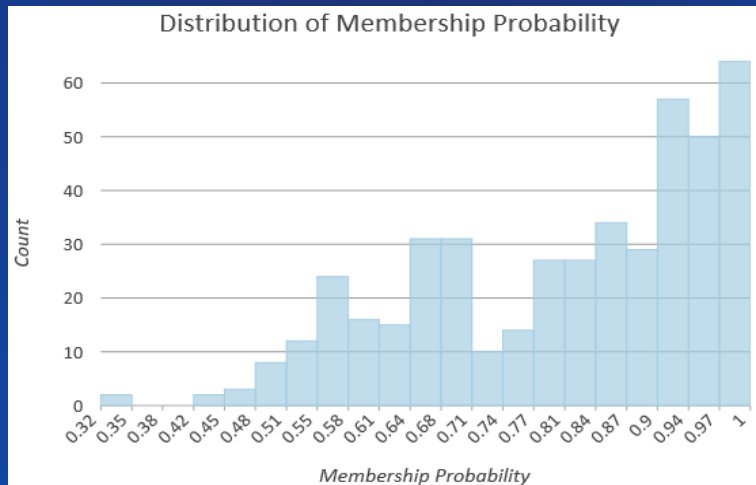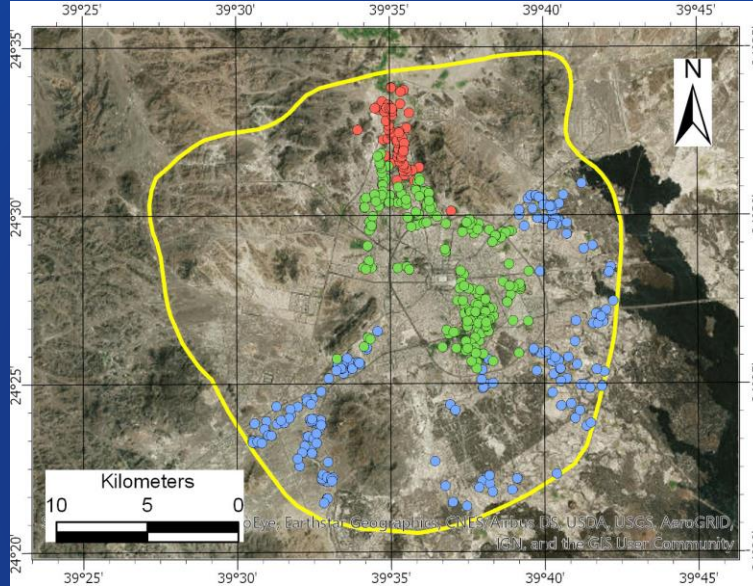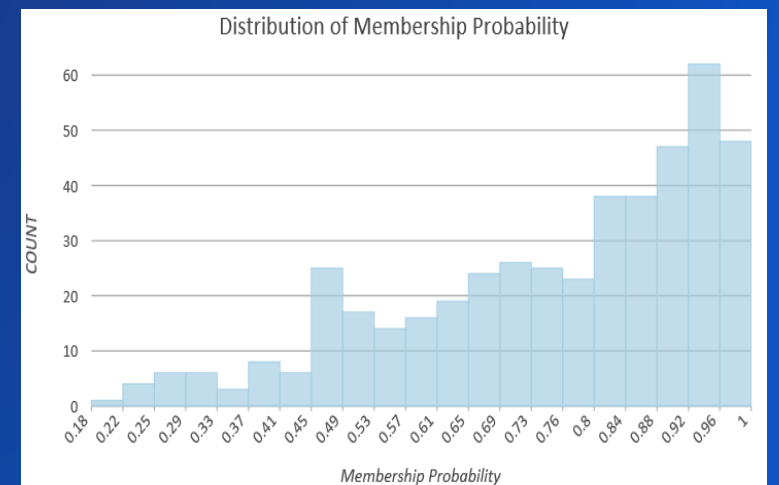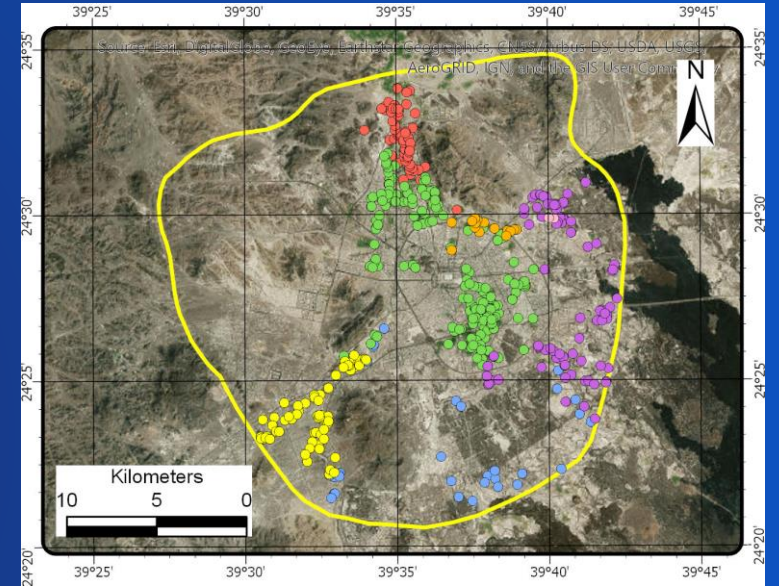
```
Variable      Mean    Std. Dev.       Min         Max        R2
     MG  184.589912 142.241319 13.000000   680.000000 0.516441
     CL 1079.949561 739.258629 88.000000  3970.000000 0.500044
     NA  542.953947 380.473044 30.000000  2520.000000 0.492381
   HCO3  271.379386 130.798051 42.000000   780.000000 0.463728
     CA  303.885965 200.843004 22.000000  1168.000000 0.417436
    SO4  982.528509 610.161313 26.000000  3996.000000 0.360552
```

# Results & Discussions (Group (C) scenarios)
Assuming the optimum No. of cluster = 7



7 clusters, no min. No. of features

7 clusters, min. No. of features = 20

7 clusters, min. No. of features = 40

# Results & Discussions (Group (C) scenarios)



7 clusters, no min. No. of features

7 clusters, min. No. of features = 20

7 clusters, min. No. of features = 40

# Results & Discussions (Group (C) scenarios)



| | 7 clusters, no min. No. of features |
|---|---|

| Variable | Mean | Std. Dev. | Min | Max | R2 |
|---|---|---|---|---|---|
| MG | 184.589912 | 142.241319 | 13.000000 | 680.000000 | 0.516441 |
| CL | 1079.949561 | 739.258629 | 88.000000 | 3970.000000 | 0.500044 |
| NA | 542.953947 | 380.473044 | 30.000000 | 2520.000000 | 0.492381 |
| HCO3 | 271.379386 | 130.798051 | 42.000000 | 780.000000 | 0.463728 |
| CA | 303.885965 | 200.843004 | 22.000000 | 1168.000000 | 0.417436 |
| SO4 | 982.528509 | 610.161313 | 26.000000 | 3996.000000 | 0.360552 |

| | 7 clusters, min. No. of features = 20 |
|---|---|

| Variable | Mean | Std. Dev. | Min | Max | R2 |
|---|---|---|---|---|---|
| MG | 184.589912 | 142.241319 | 13.000000 | 680.000000 | 0.516441 |
| CL | 1079.949561 | 739.258629 | 88.000000 | 3970.000000 | 0.500044 |
| NA | 542.953947 | 380.473044 | 30.000000 | 2520.000000 | 0.492381 |
| HCO3 | 271.379386 | 130.798051 | 42.000000 | 780.000000 | 0.463728 |
| CA | 303.885965 | 200.843004 | 22.000000 | 1168.000000 | 0.417436 |
| SO4 | 982.528509 | 610.161313 | 26.000000 | 3996.000000 | 0.360552 |

| | 7 clusters, min. No. of features = 40 |
|---|---|

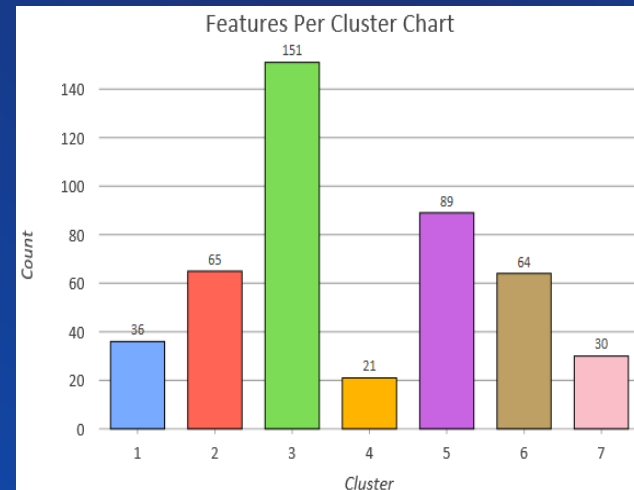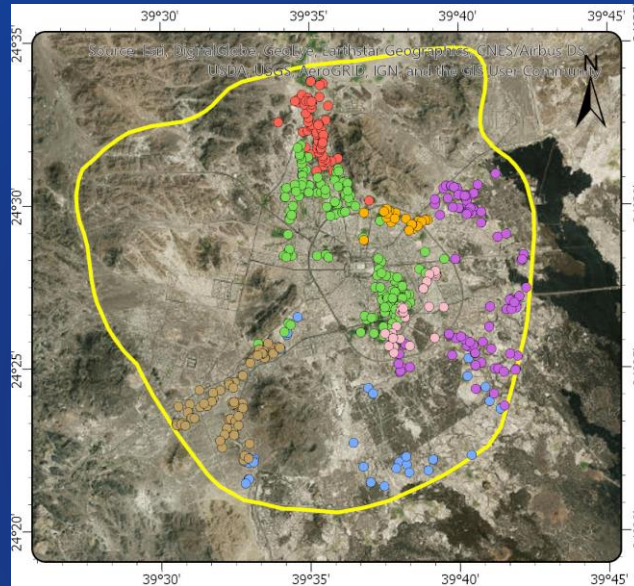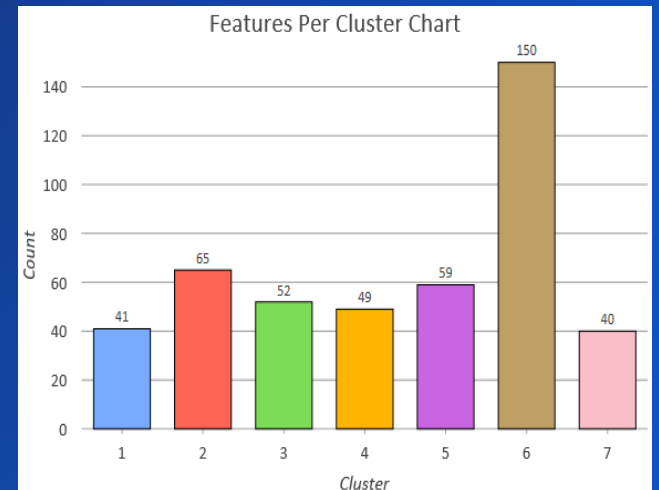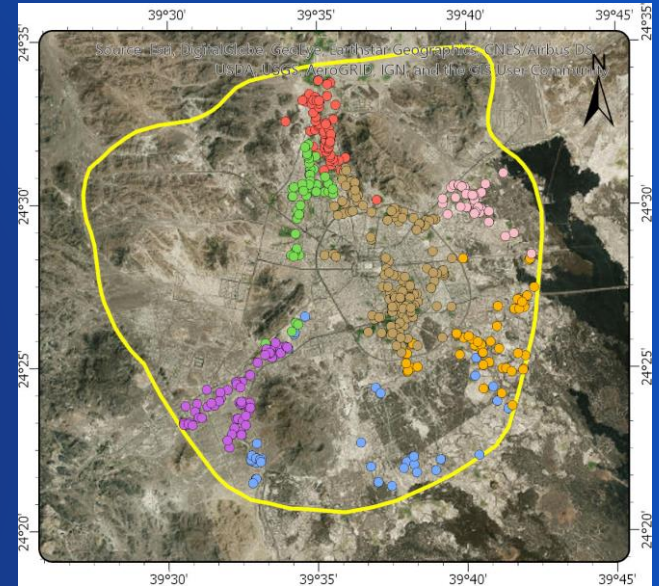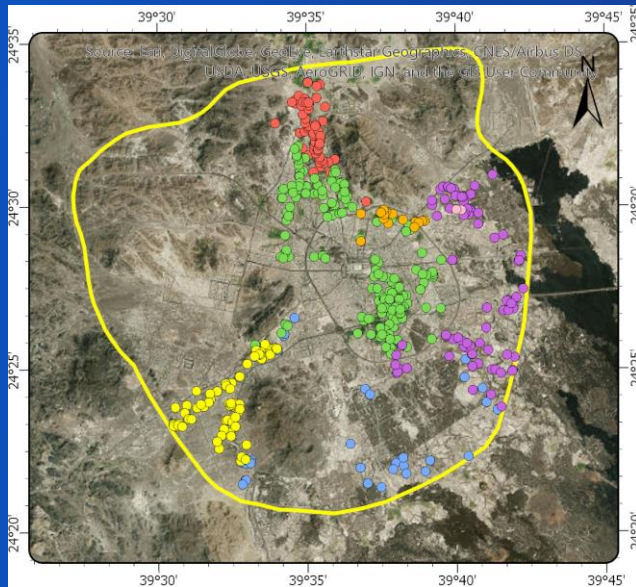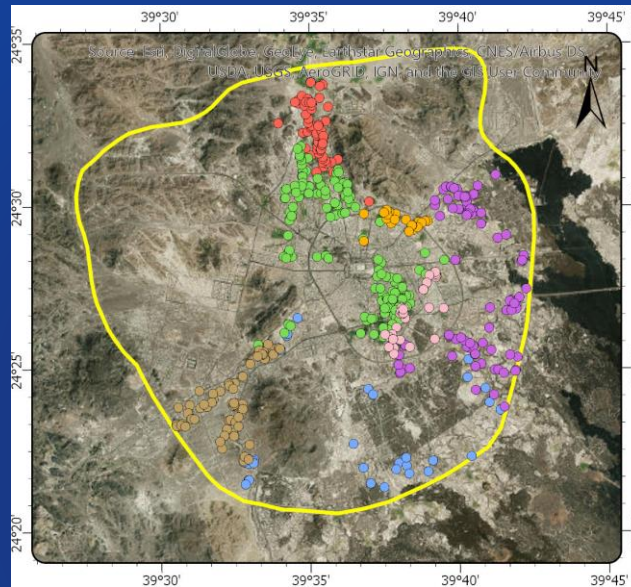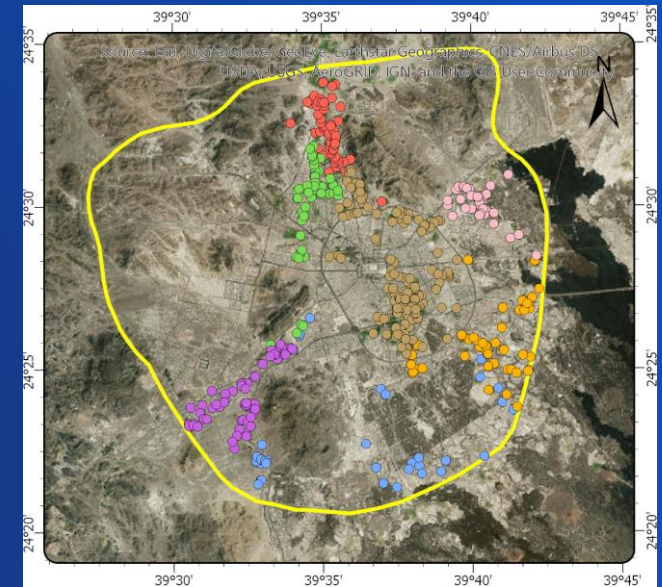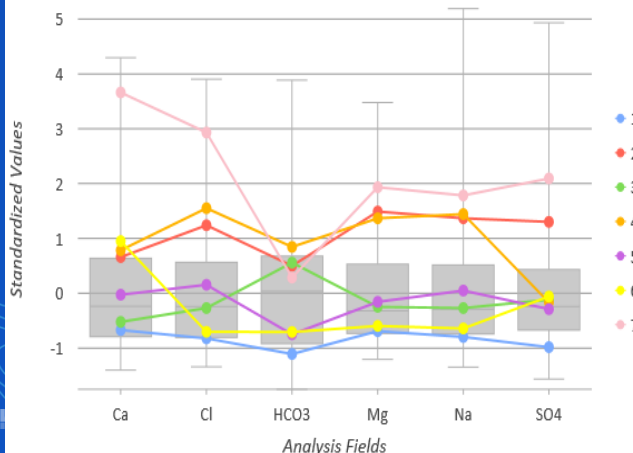| Variable | Mean | Std. Dev. | Min | Max | R2 |
|---|---|---|---|---|---|
| HCO3 | 271.379386 | 130.798051 | 42.000000 | 780.000000 | 0.501893 |
| MG | 184.589912 | 142.241319 | 13.000000 | 680.000000 | 0.436772 |
| NA | 542.953947 | 380.473044 | 30.000000 | 2520.000000 | 0.410843 |
| CL | 1079.949561 | 739.258629 | 88.000000 | 3970.000000 | 0.395729 |
| SO4 | 982.528509 | 610.161313 | 26.000000 | 3996.000000 | 0.354156 |
| CA | 303.885965 | 200.843004 | 22.000000 | 1168.000000 | 0.294929 |

# Results & Discussions (Group (D) scenarios)

Specifying the margins (boundaries) of the No. of features per cluster

Min. No. of features/cluster = 25
Max. No. of features/cluster = 150

Min. No. of features/cluster = 50
Max. No. of features/cluster = 100

# Results & Discussions (Group (D) scenarios)

Min. No. of features/cluster = 25
Max. No. of features/cluster = 150

Min. No. of features/cluster = 50
Max. No. of features/cluster = 100

# Results & Discussions (Group (D) scenarios)

Min. No. of features/cluster = 25
Max. No. of features/cluster = 150

Min. No. of features/cluster = 50
Max. No. of features/cluster = 100



| Variable | Mean | Std. Dev. | Min | Max | R2 |
|---|---|---|---|---|---|
| HCO3 | 271.379386 | 130.798051 | 42.000000 | 780.000000 | 0.494480 |
| MG | 184.589912 | 142.241319 | 13.000000 | 680.000000 | 0.427112 |
| NA | 542.953947 | 380.473044 | 30.000000 | 2520.000000 | 0.405020 |
| CL | 1079.949561 | 739.258629 | 88.000000 | 3970.000000 | 0.390493 |
| SO4 | 982.528509 | 610.161313 | 26.000000 | 3996.000000 | 0.342758 |
| CA | 303.885965 | 200.843004 | 22.000000 | 1168.000000 | 0.174965 |



| Variable | Mean | Std. Dev. | Min | Max | R2 |
|---|---|---|---|---|---|
| MG | 184.589912 | 142.241319 | 13.000000 | 680.000000 | 0.426960 |
| HCO3 | 271.379386 | 130.798051 | 42.000000 | 780.000000 | 0.412836 |
| NA | 542.953947 | 380.473044 | 30.000000 | 2520.000000 | 0.398952 |
| CL | 1079.949561 | 739.258629 | 88.000000 | 3970.000000 | 0.388360 |
| SO4 | 982.528509 | 610.161313 | 26.000000 | 3996.000000 | 0.293898 |
| CA | 303.885965 | 200.843004 | 22.000000 | 1168.000000 | 0.278877 |

# Results & Discussions (more spatial control using spatial weights)



- **Weighted Optimization**

  $w_1$(attribute similarity) + $w_2$(geometric centroids)

  $w_1 + w_2 = 1$

  iterate until contiguity constraint is satisfied

  bisection method

  $w_2$ is weight for centroids, $w_1 = 1 - w_2$

  start with 0.0 and 1.0

  then move to 0.50 - check contiguity

  if contiguous, then to midpoint to the left of 0.50

  if not contiguous, then to midpoint to the right of 0.50

  etc... until contiguous with the highest bSS/tSS ratio

20

**Cluster Size Constraints**
None

**Number of Clusters**

ⓘ Spatial Constraints
Trimmed Delaunay triangulation

Trimmed Delaunay triangulation
Get spatial weights from file

Output Table for Evaluating Number of Clusters

# Results & Discussions

- Higher Goodness of fit index is better
- B = between-cluster sum of square error (SSE) – need to be maximized
- W = within-cluster sum of square error (SSE)
- K = the number of clusters
- N = the number of features (observation)

$$Goodness\ of\ fit\ idex = \frac{B/(k-1)}{W/(n-k)}$$

# Conclusion

- **SCMC method in ArcGIS Pro found to be powerful tool for spatial clustering with many options and functionalities.**

- **SCMC method as many Clustering methods needs full understanding of the data used.**

- **One of the best scenarios is sub-dividing the city based on GW quality into three zones which are; Upper zone with good quality, city center zone with moderate quality, and lower (downstream) zone with low quality.**

- **The results of this study will be beneficial not only for the farmers but also for the local government, environmental agencies and investors in agriculture.**

# References and Resources

- **Spatially Constrained Multivariate Clustering**
    - https://pro.arcgis.com/en/pro-app/tool-reference/spatial-statistics/spatially-constrained-multivariate-clustering.htm

- **How Spatially Constrained Multivariate Clustering works**
    - https://pro.arcgis.com/en/pro-app/tool-reference/spatial-statistics/how-spatially-constrained-multivariate-clustering-works.htm

- **Determining Number of Clusters in One Picture**
    - https://www.datasciencecentral.com/profiles/blogs/determining-number-of-clusters-in-one-picture

- **Assunção, R. M., Neves, M. C., Câmara, G., & da Costa Freitas, C. (2006). Efficient regionalization techniques for socio-economic geographical units using minimum spanning trees. International Journal of Geographical Information Science, 20(7), 797-811.**

- **Duque, J. C., Ramos, R., & Suriñach, J. (2007). Supervised regionalization methods: A survey. *International Regional Science Review*, *30*(3), 195-220.**

- **_Luc Anselin, (2017),_ Cluster Analysis (3) *Spatially Constrained Clustering Methods,***
    - https://geodacenter.github.io/workbook/8_spatial_clusters/lab8.html

- **Kassambara, A. (2017). *Practical guide to cluster analysis in R: unsupervised machine learning* (Vol. 1). STHDA.**

- **Fischer, M. M., & Getis, A. (Eds.). (2009). *Handbook of applied spatial analysis: software tools, methods and applications*. Springer Science & Business Media.**

# Thank you