



# Streamlining Principal Component and Cluster Analysis to Support Conceptual Site Models at Complex Remediation Sites

Tori Ward

May 14, 2025

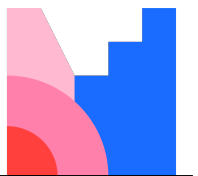
# Introduction

---

# Menti Question ([www.menti.com](http://www.menti.com)):

How familiar are you with PCA and Cluster analysis?

**6927 2321**

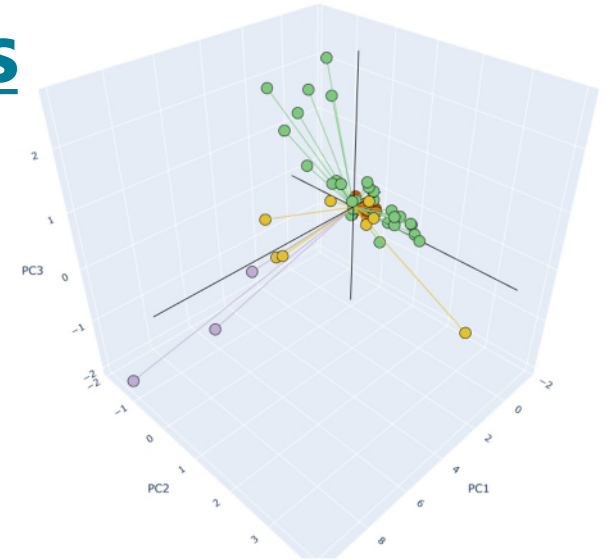


# Intro to PCA and HCA

## Unsupervised Machine Learning Methods

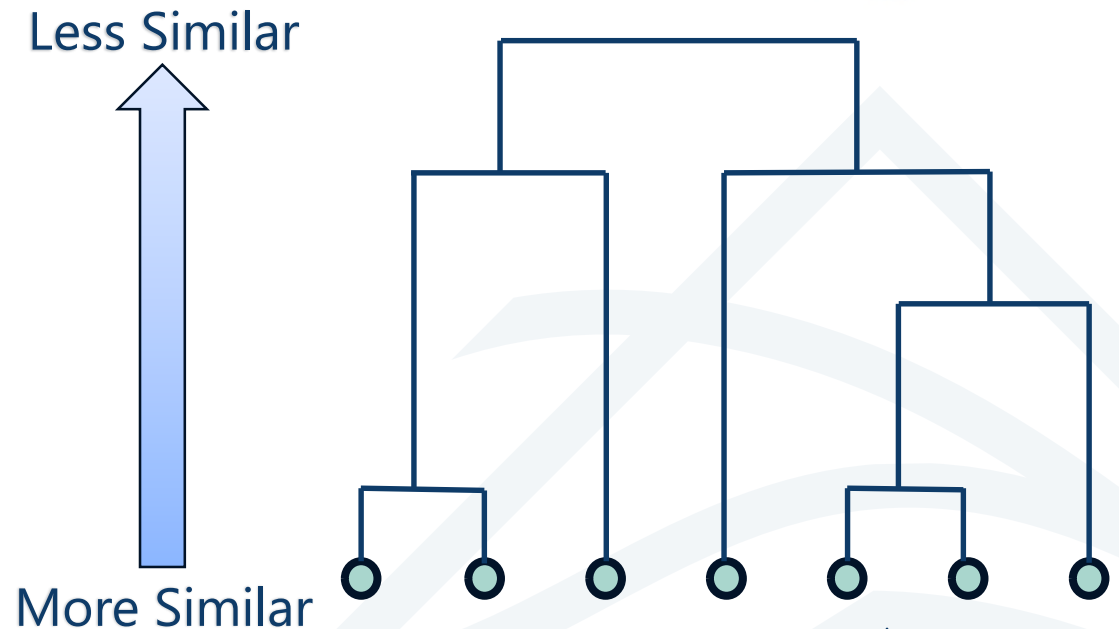
- **Principal Component Analysis (PCA)**

- Reduces many variables to a few components that explain most of the variability in the data
- **Components** are linear combinations of variables, ordered by decreasing importance



- **Hierarchical Cluster Analysis (HCA)\***

- Starts by placing each observation in its own **cluster** (Bottom-up approach)
- Applies a distance measure to merge clusters that are most similar



• *\*HCA is just one example of a cluster analysis*

# Why use PCA and Cluster Analysis?



**Complex Geology**



**Multiple Suspected Sources**

*\*not a real Site*

## **Menti Question ([www.menti.com](http://www.menti.com)):**

What are other ways you are aware of PCA and cluster analysis being used in our industry?

**6927 2321**

# Methodology and Workflow

---

# python™ Tools for PCA and HCA

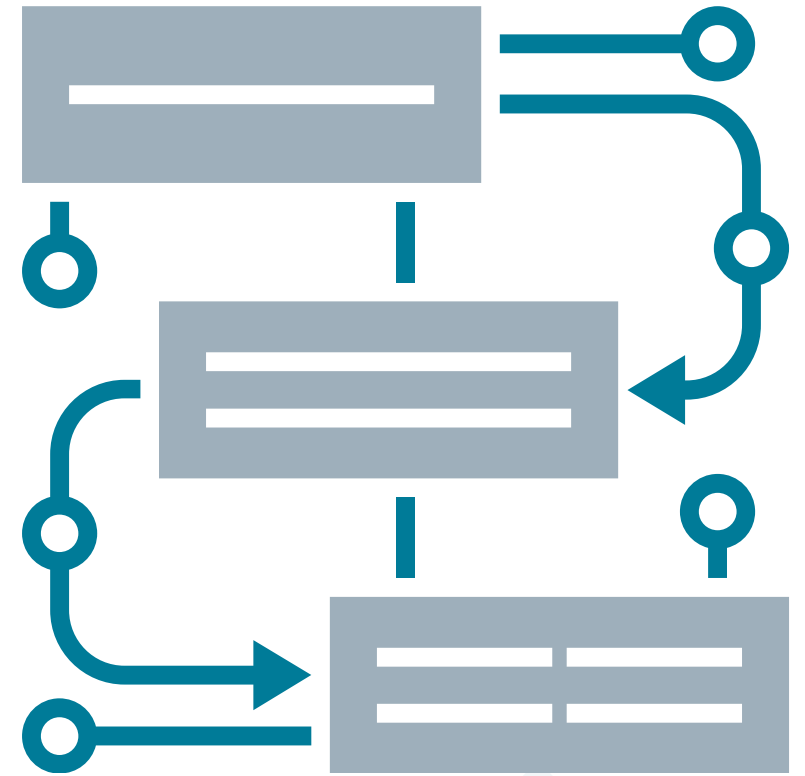


[This Photo](#) by Unknown Author is licensed under [CC BY](#)

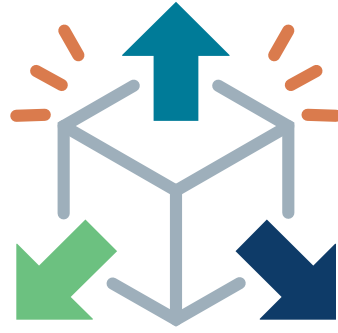
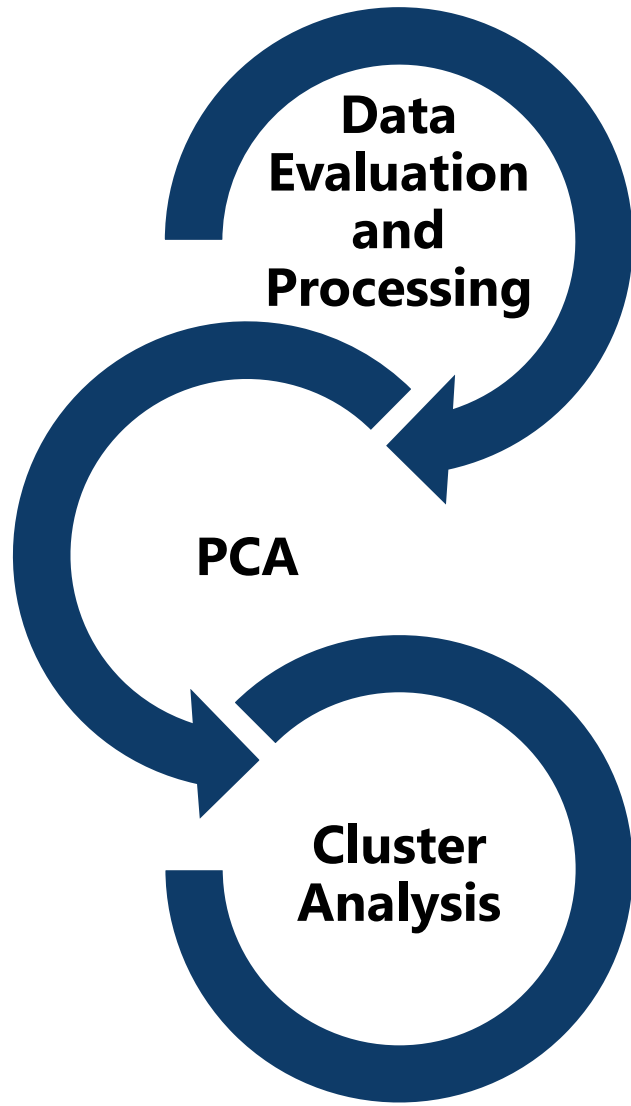


# Standardized Methodology

- ▶ Basic statistics
- ▶ Correlation matrix
- ▶ PCA
  - Data transformations
  - Review components, explained variance, and loadings
- ▶ Visualize PCA
  - 3D charts of PCA factor loadings
- ▶ Clustering
  - Selecting an appropriate number of clusters
    - Dendrograms
    - Elbow and Silhouette Scores
    - Professional Judgement when reviewing output
  - Applied a color-scale based on cluster labels to 2D and 3D PCA charts to assess corroboration between PCA and agglomerative clustering results



# Streamlined Workflow



**Project-specific customizations**

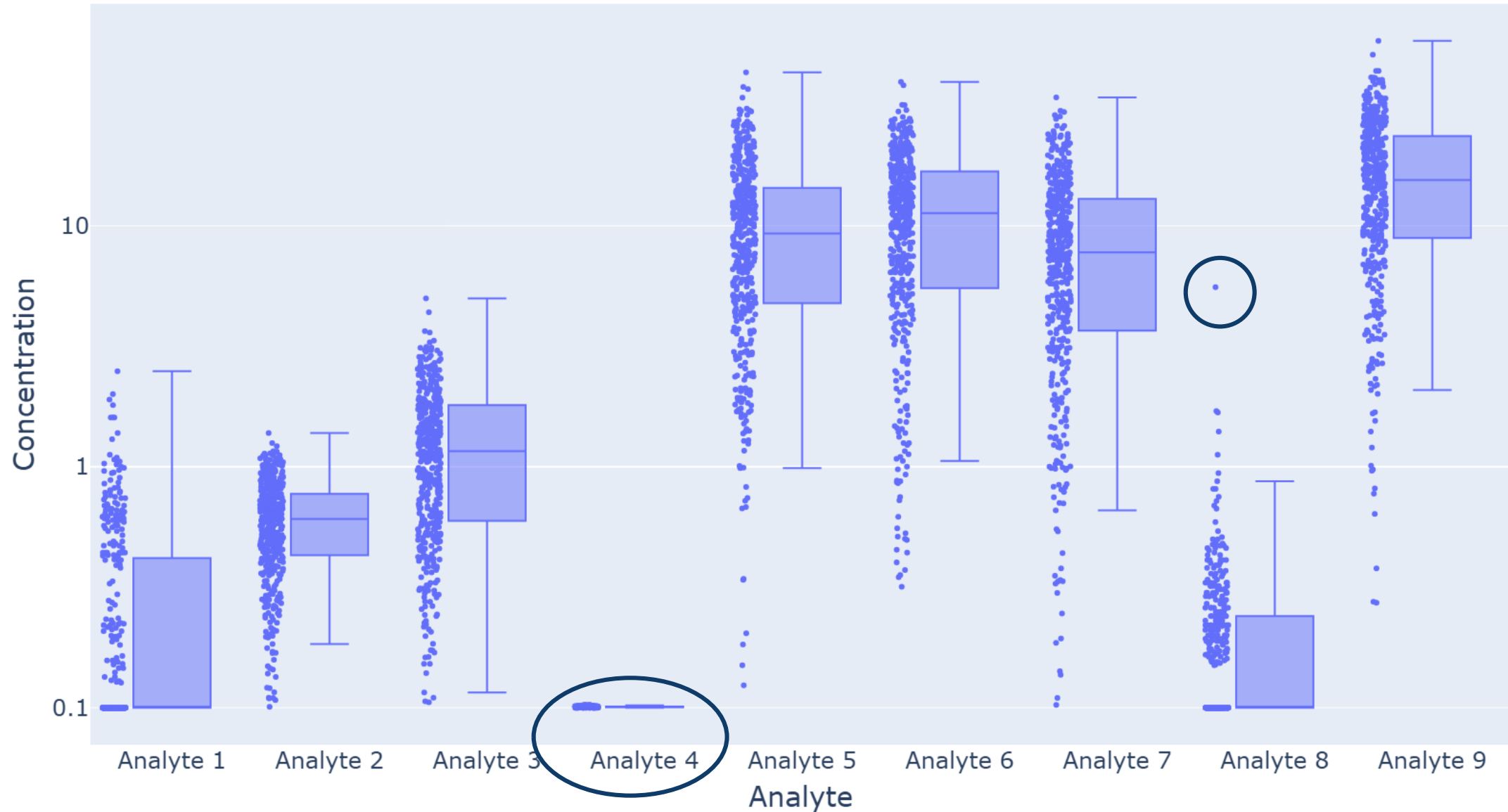
**Generalized functions across all 3 scripts are stored in a utilities file**

**Code maintained on GitHub** 

# Example Output Visualizations

*\*Graphics and tables presented do not summarize project data to preserve client confidentiality*

# Methodology – Basic Stats



# Methodology – Basic Stats

<b>Analyte</b>	<b>Number of Detections</b>	<b>Number of Analyses</b>	<b>Frequency of Detection</b>
Analyte 6	404	442	91.40%
Analyte 7	350	442	79.19%
Analyte 5	311	442	70.36%
Analyte 3	281	436	64.45%
Analyte 2	236	442	53.39%
Analyte 9	211	415	50.84%
Analyte 8	185	410	45.12%
Analyte 4	163	408	39.95%
Analyte 1	155	415	37.35%

## **Menti Question ([www.menti.com](http://www.menti.com)):**

Which methods have you used to impute non-detect results (multiple selections possible)?

**6927 2321**



# Methodology – Correlation Matrix

Analyte	Analyte 1	Analyte 2	Analyte 3	Analyte 4	Analyte 5	Analyte 6
Analyte 1	1.00					
Analyte 2	0.33	1.00				
Analyte 3	<b>0.70</b>	0.39	1.00			
Analyte 4	0.34	0.33	0.40	1.00		
Analyte 5	<b>0.78</b>	0.66	<b>0.73</b>	0.27	1.00	
Analyte 6	0.42	0.39	0.64	0.21	0.48	1.00

**Key Point:** High correlation between certain analytes suggests they can be aggregated into a smaller number of components

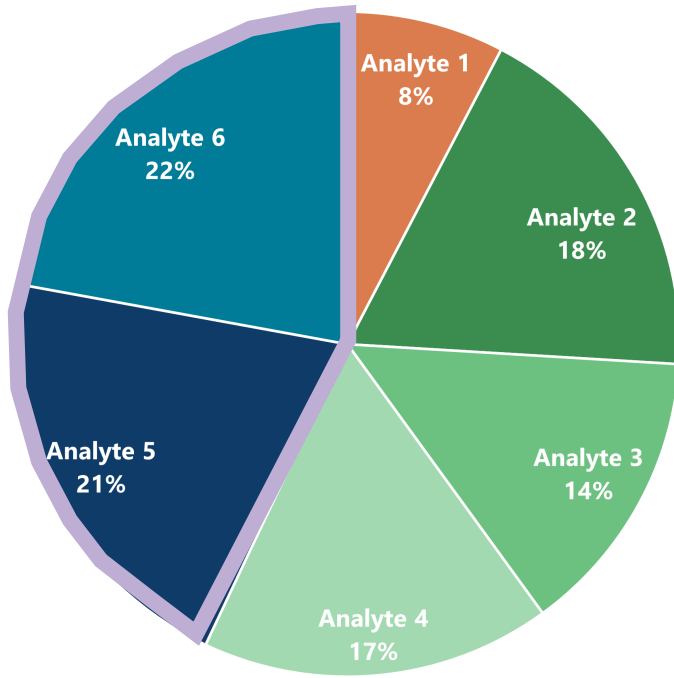
# Methodology – PCA Component Weights

	Analyte 1	Analyte 2	Analyte 3	Analyte 4	Analyte 5	Analyte 6	Explained Variance
<b>PC1</b>	<b>0.18</b>	<b>0.43</b>	<b>0.33</b>	<b>0.40</b>	<b>0.49</b>	<b>0.52</b>	<b>58.0%</b>
<b>PC2</b>	<b>0.95</b>	<b>-0.11</b>	<b>-0.13</b>	<b>0.14</b>	<b>-0.19</b>	<b>-0.10</b>	<b>15.6%</b>
<b>PC3</b>	<b>0.12</b>	<b>-0.41</b>	<b>0.85</b>	<b>-0.32</b>	<b>0.04</b>	<b>-0.04</b>	<b>13.3%</b>
<b>PC4</b>	<b>0.20</b>	<b>0.52</b>	<b>-0.09</b>	<b>-0.77</b>	<b>0.29</b>	<b>-0.12</b>	<b>9.4%</b>
<b>PC5</b>	<b>-0.05</b>	<b>0.61</b>	<b>0.38</b>	<b>0.15</b>	<b>-0.64</b>	<b>-0.24</b>	<b>2.9%</b>
<b>PC6</b>	<b>0.02</b>	<b>-0.05</b>	<b>-0.09</b>	<b>-0.33</b>	<b>-0.48</b>	<b>0.80</b>	<b>0.9%</b>

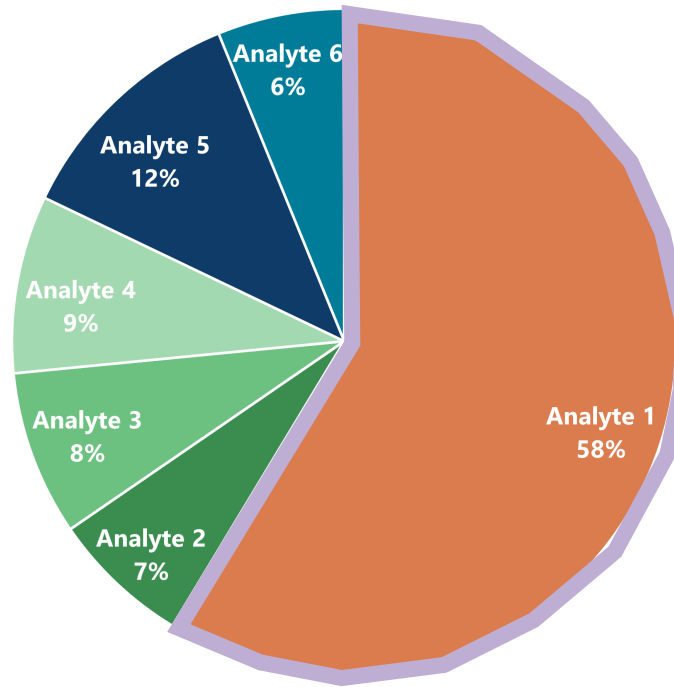
**Key Point:** Analytes 1, 3, 5, and 6 explain most of the variability in the original data set.

# Methodology – PCA Component Weights Visualization

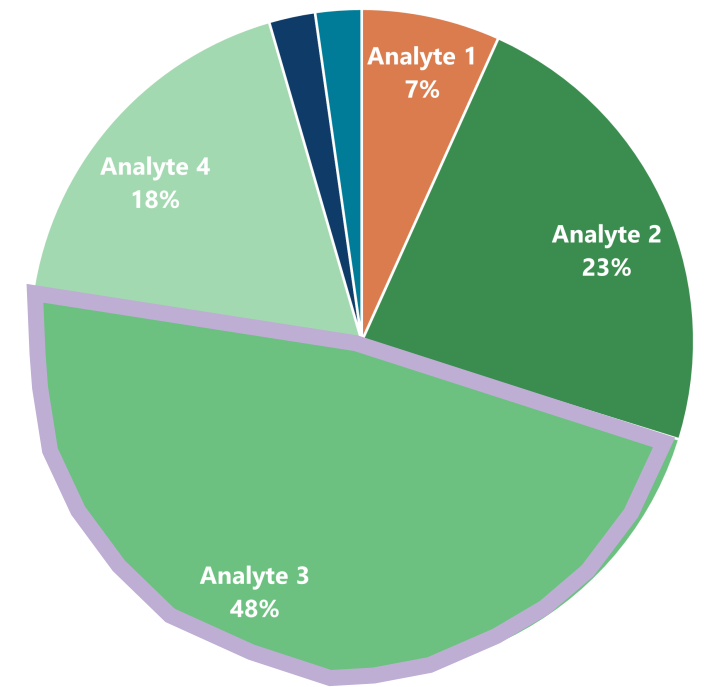
## Principal Component 1



## Principal Component 2



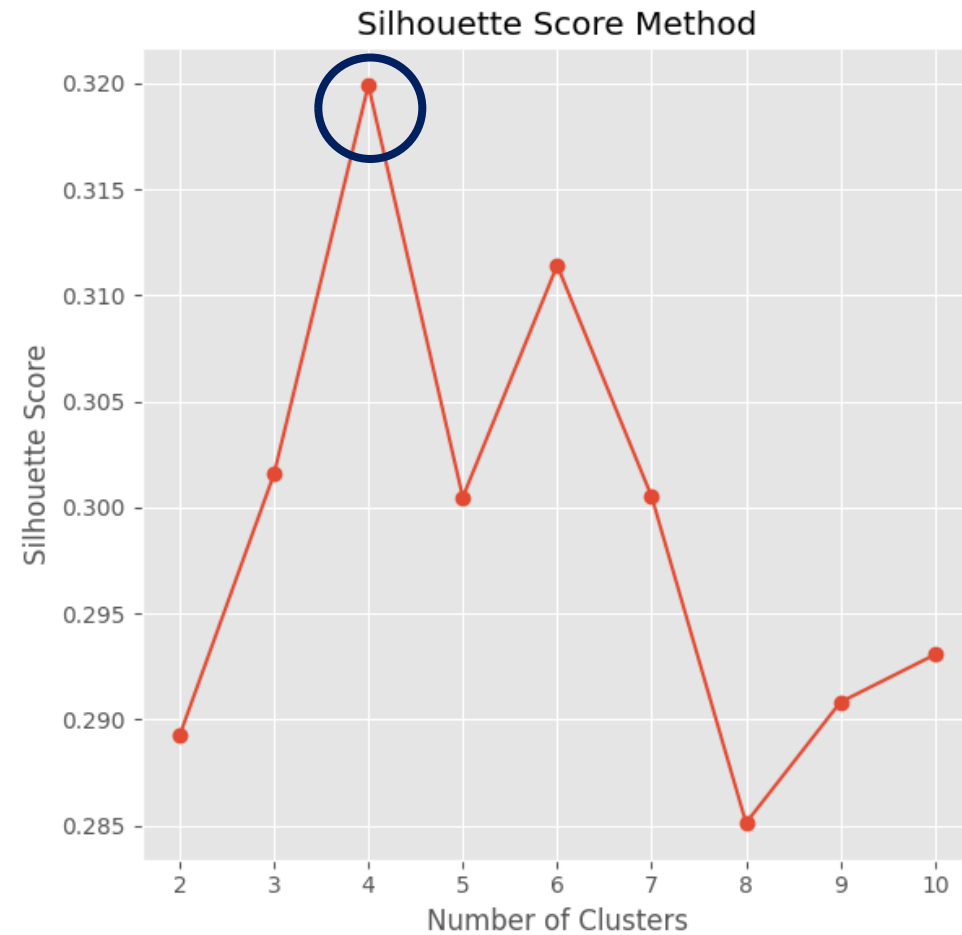
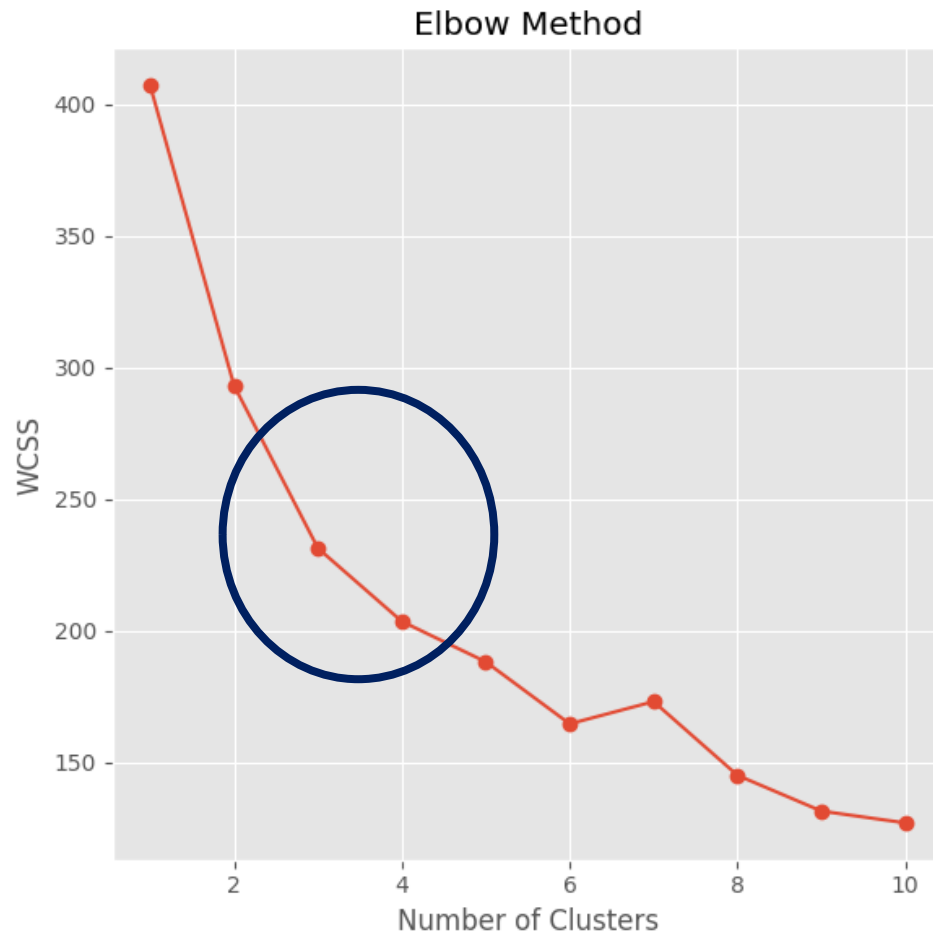
## Principal Component 3



## Key Points:

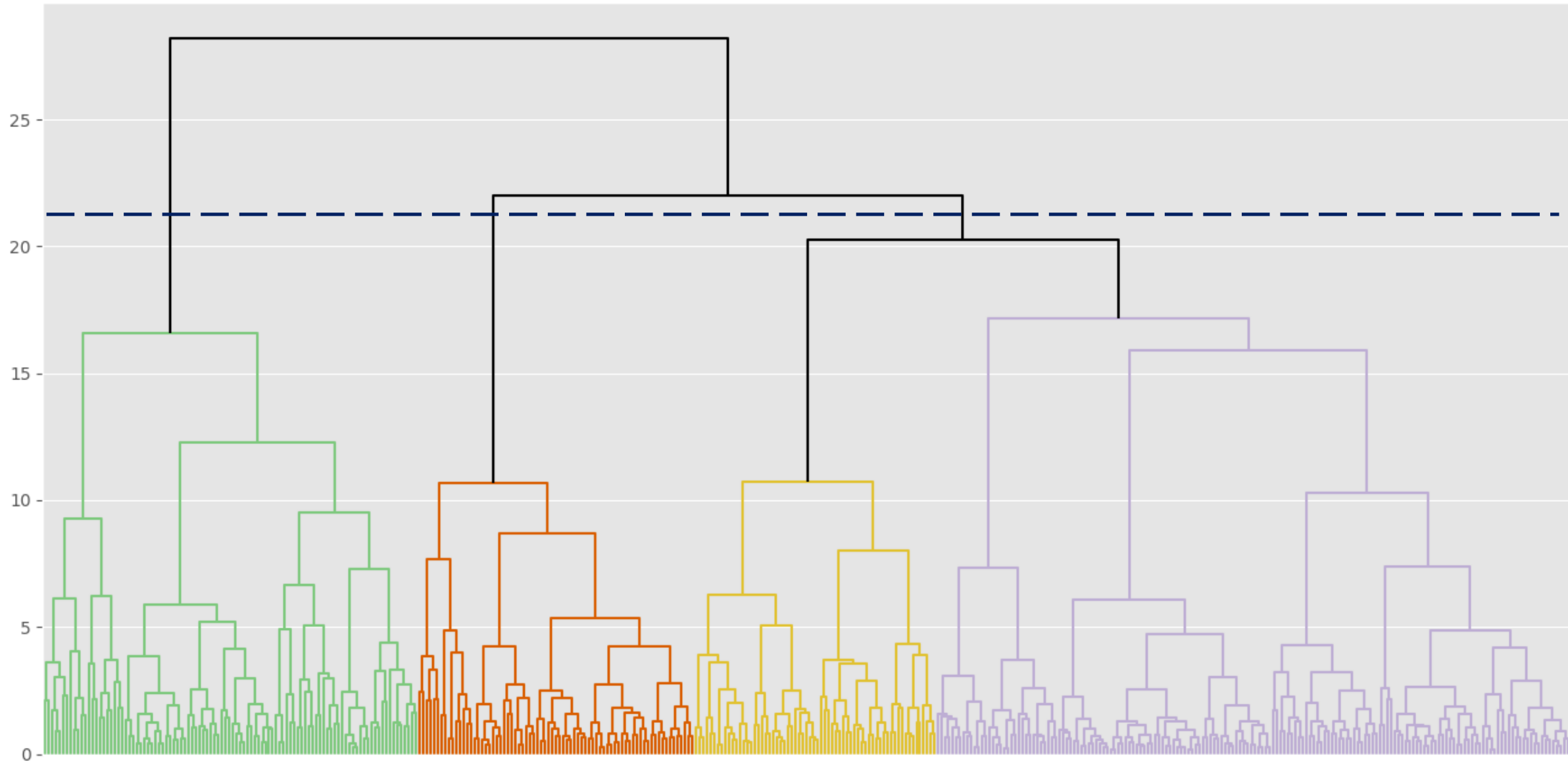
- Principal Component 1 is strongly tied to Analytes 5 and 6
- Principal Component 2 is strongly related to Analyte 1
- Principal Component 3 is strongly related to Analyte 3

# Results – Cluster Analysis Elbow and Silhouette Charts

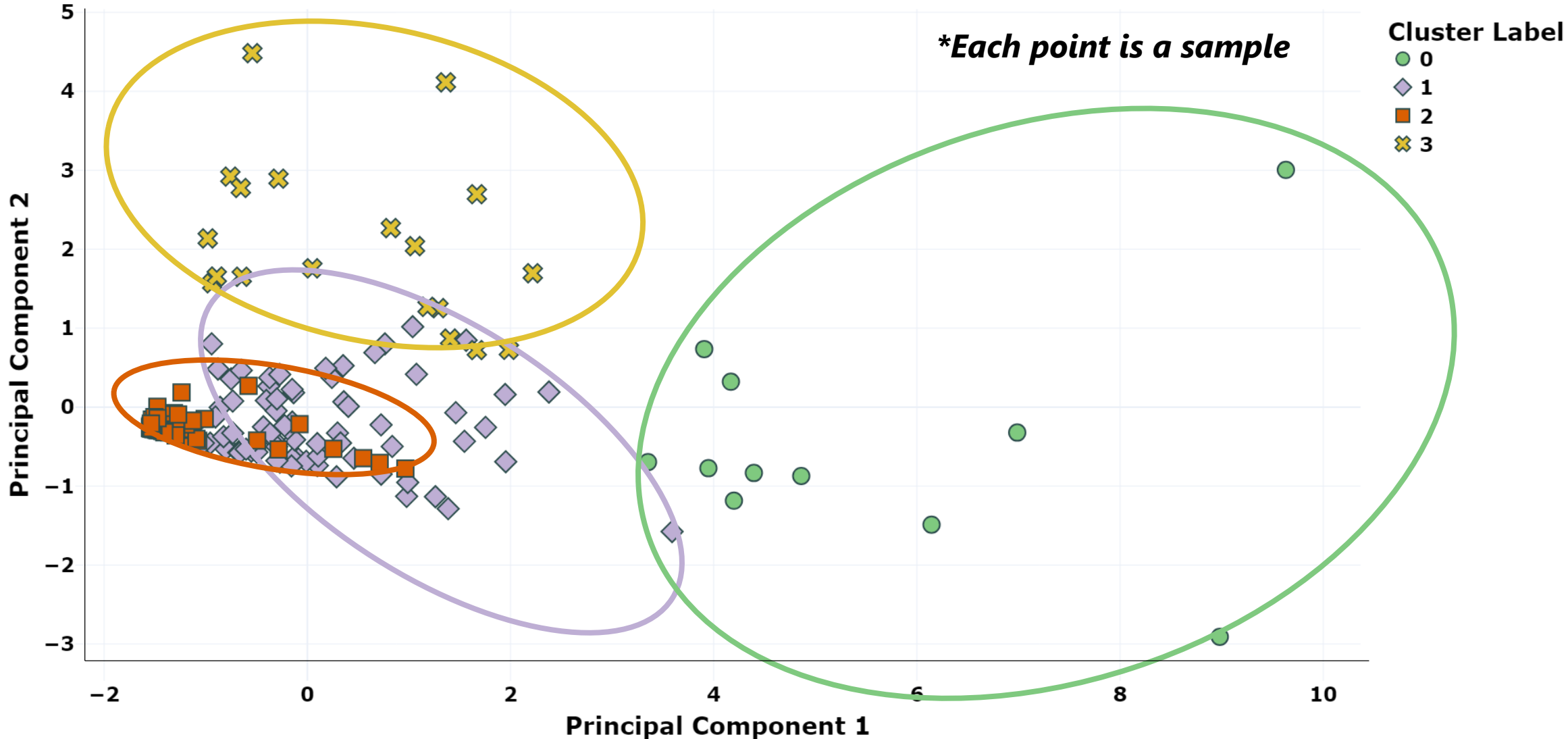


**Key Point:** 3 to 4 clusters is reasonable

# Results - HCA Dendrogram



# Results – PCA Scatter with Cluster Labels



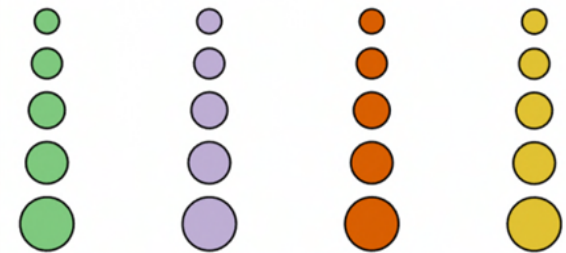
# Results – Geospatial Evaluation of Cluster Labels



## Key Points:

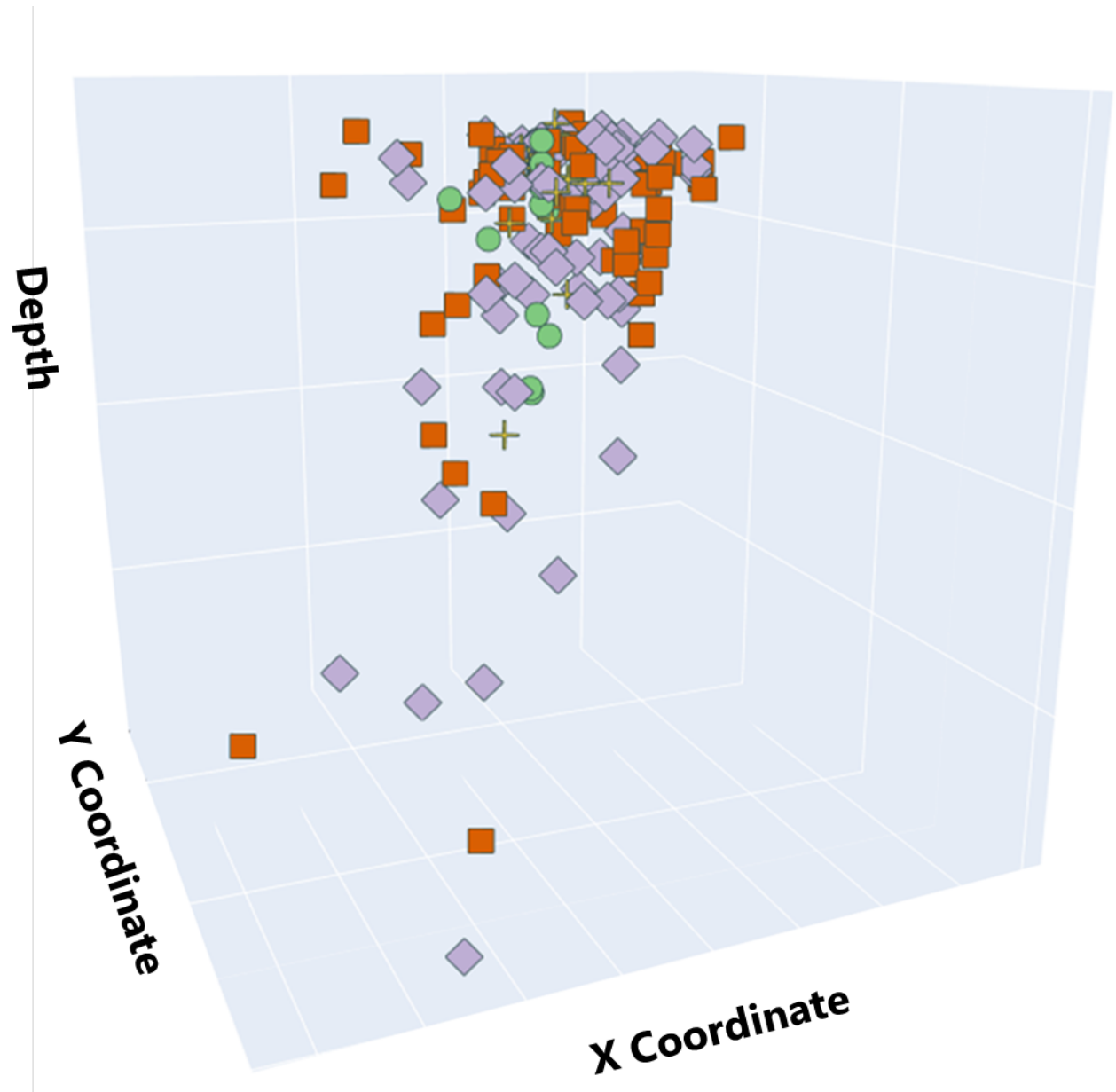
- Locations associated with clusters 0, 1, and 2 are generally located near each of the source areas
- Contaminant concentrations are generally highest in Cluster 1 and lowest in Cluster 0

Cluster 0 Cluster 1 Cluster 2 Cluster 3



*\* Not a real Site*

# X,Y,Z-Coordinates with Cluster Labels



## Cluster Label

- 0
- ◆ 1
- 2
- + 3

*\*Each point is a sample*

# Relevance

- Focus on a few key variables in a large data set
- Objective, quantitative line of evidence
- Can use results to target site remediation and planning





Tori Ward



[vward@woodardcurran.com](mailto:vward@woodardcurran.com)

Questions?