



XPLORATORY



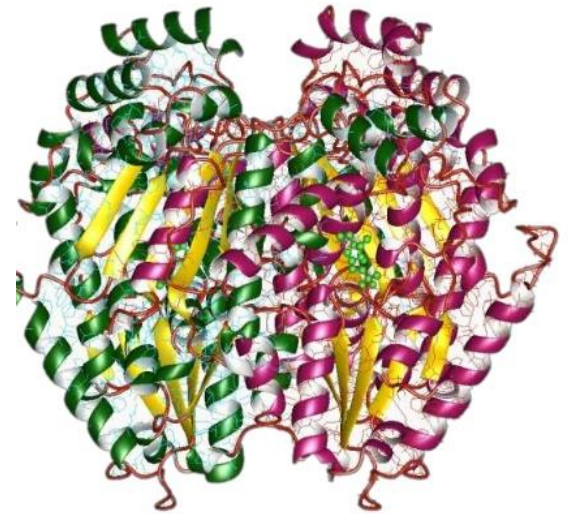
ATA



NALYSIS

@luminousmen.com

ENZYMES



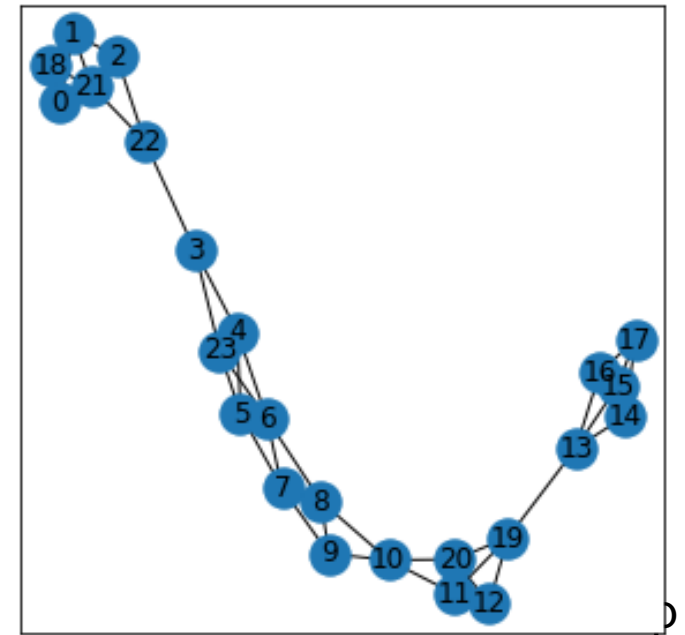
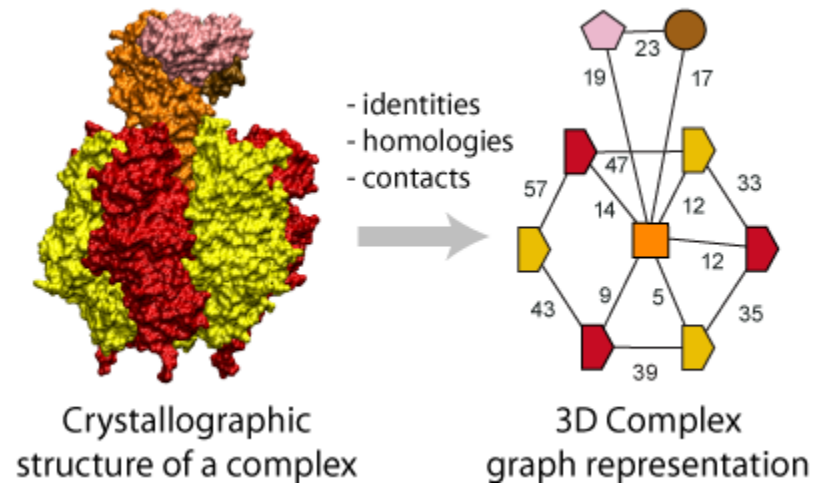
Peter Hunyadi & Nikola Vinko

ENZYMES Overview

- 600 enzymes
- Graph representation of an enzyme
- 6 EC enzyme labels: reflect catalyzed chemical reaction

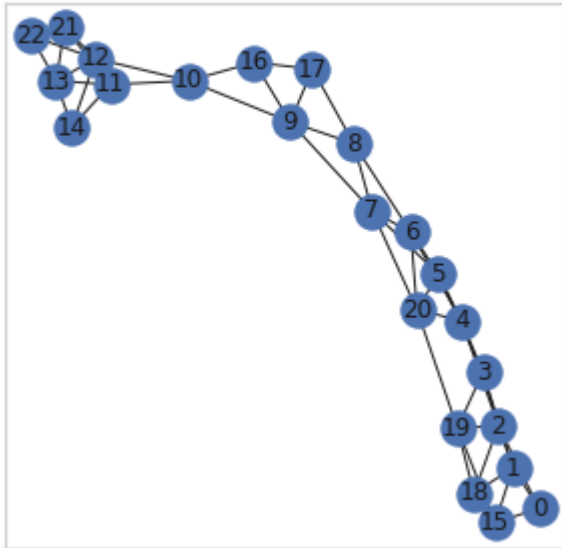
Example graph

- Node represent secondary structure elements (SSE)
 - Helix, sheet or turn
- Edge connects two SSE that are neighbors along the amino acid sequence or are one of three nearest neighbors in space

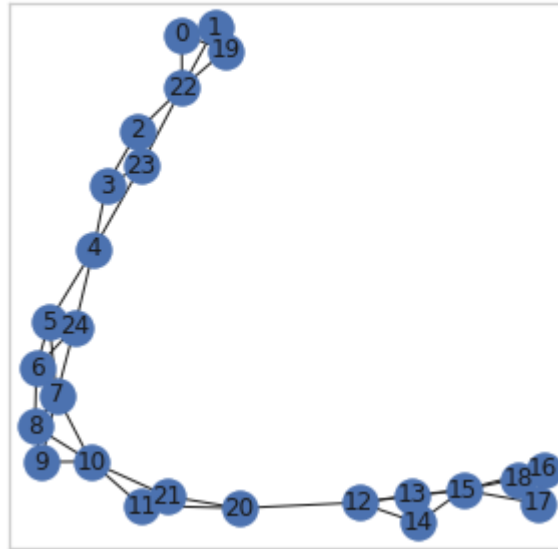


Examples of representative enzymes

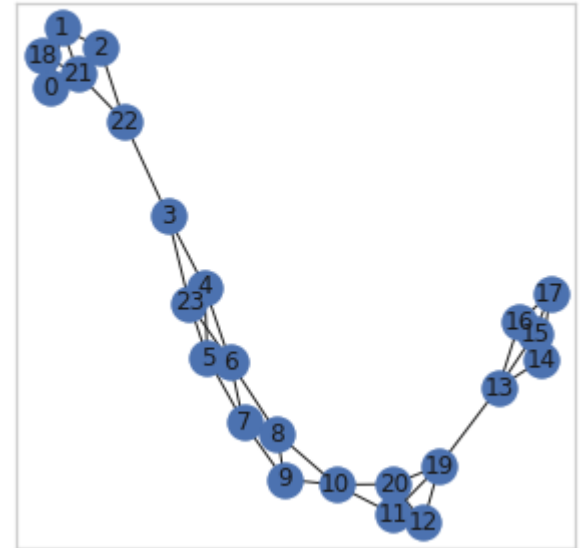
EC1



EC2

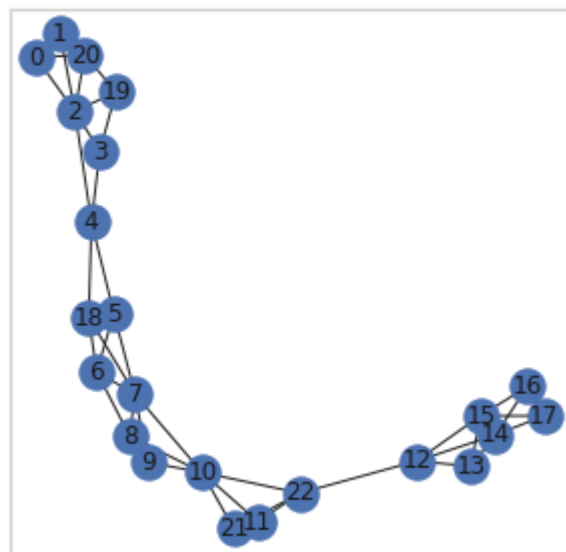


EC3

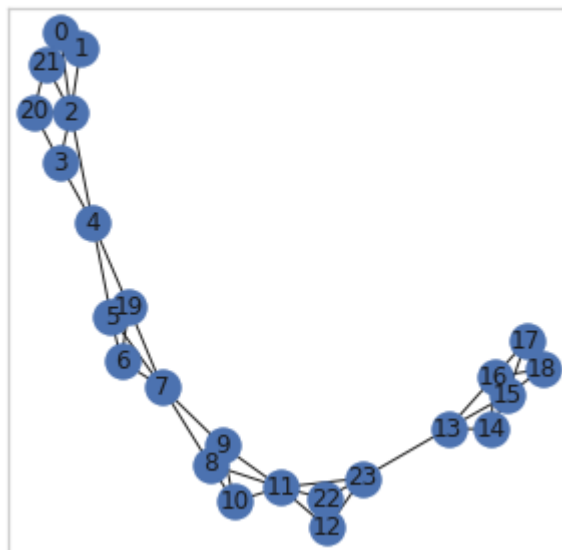


Examples of representative enzymes

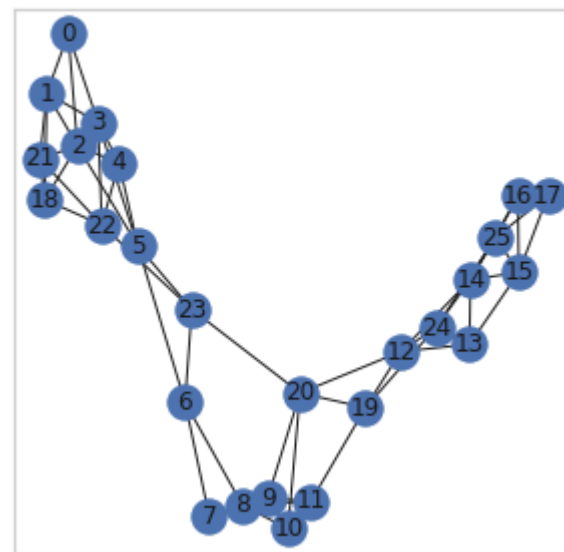
EC4



EC5

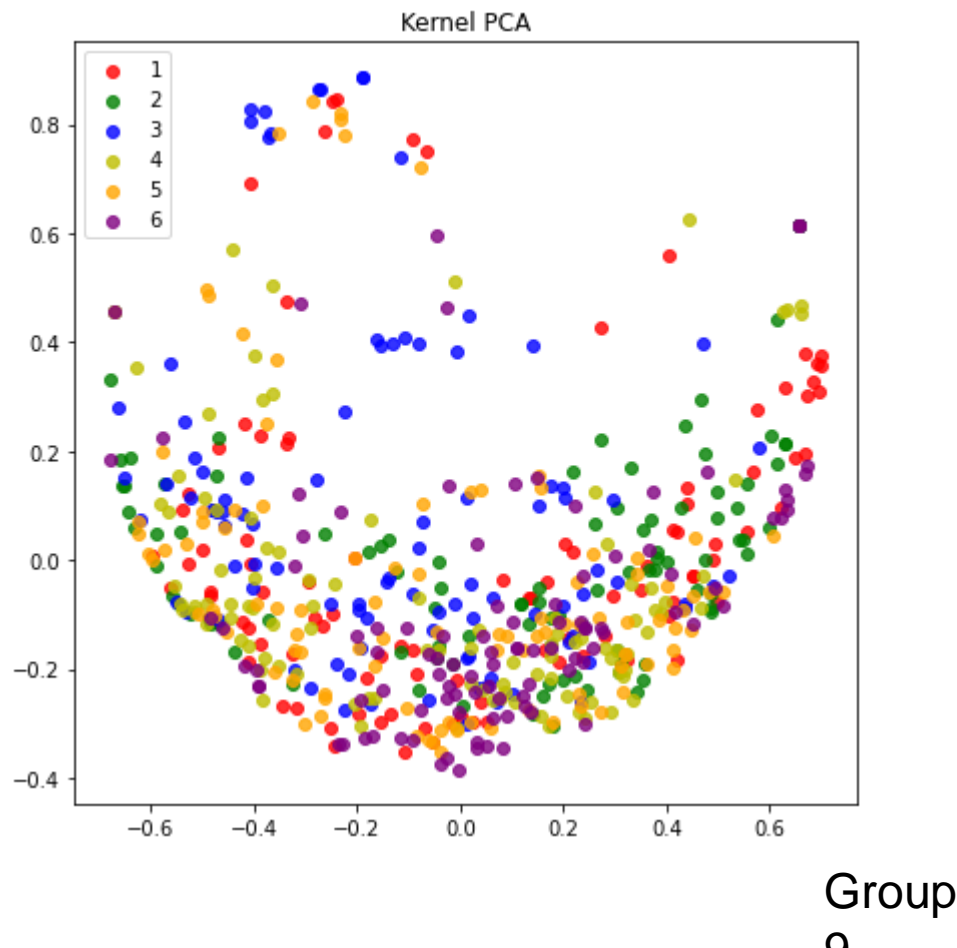


EC6

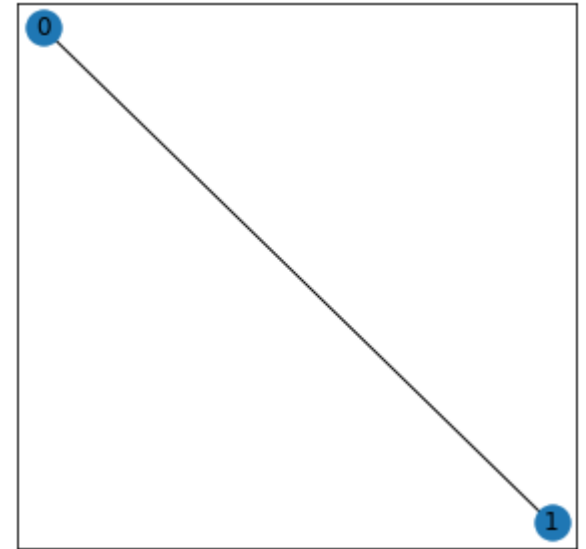
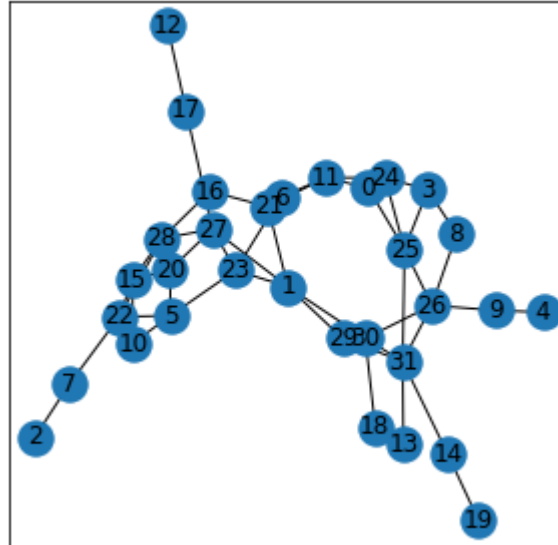
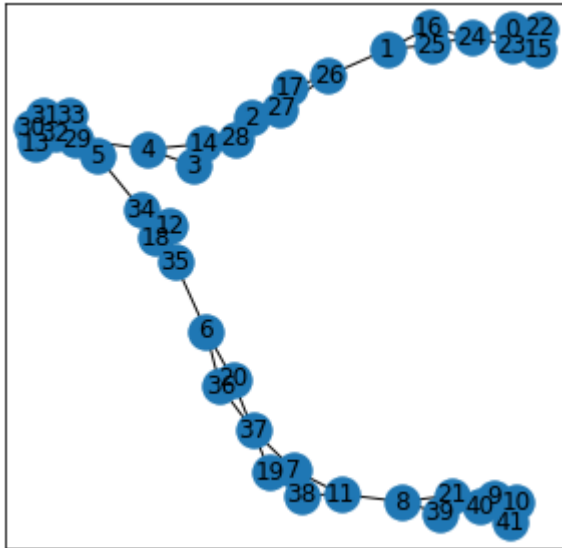


Dimensionality Reduction on unprocessed data

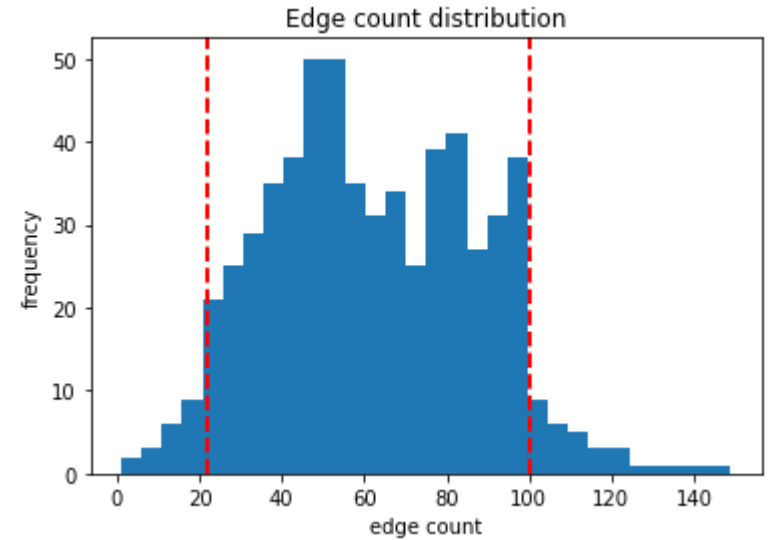
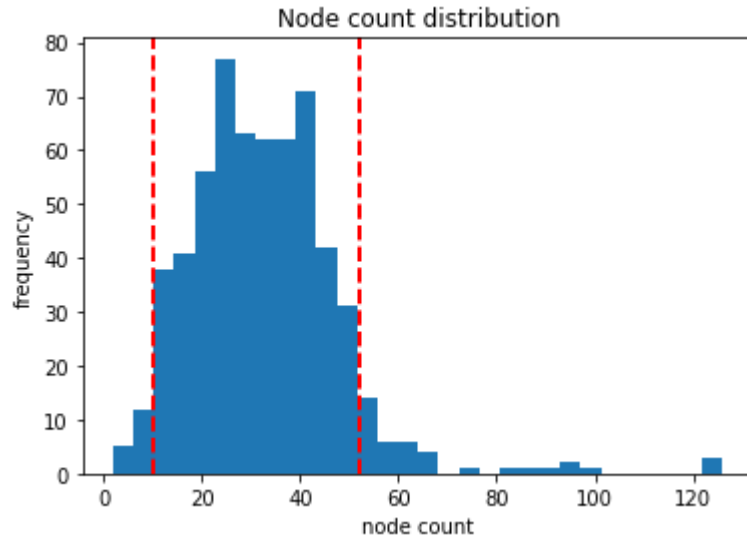
- We see some outlier points which are far separated from the main cluster



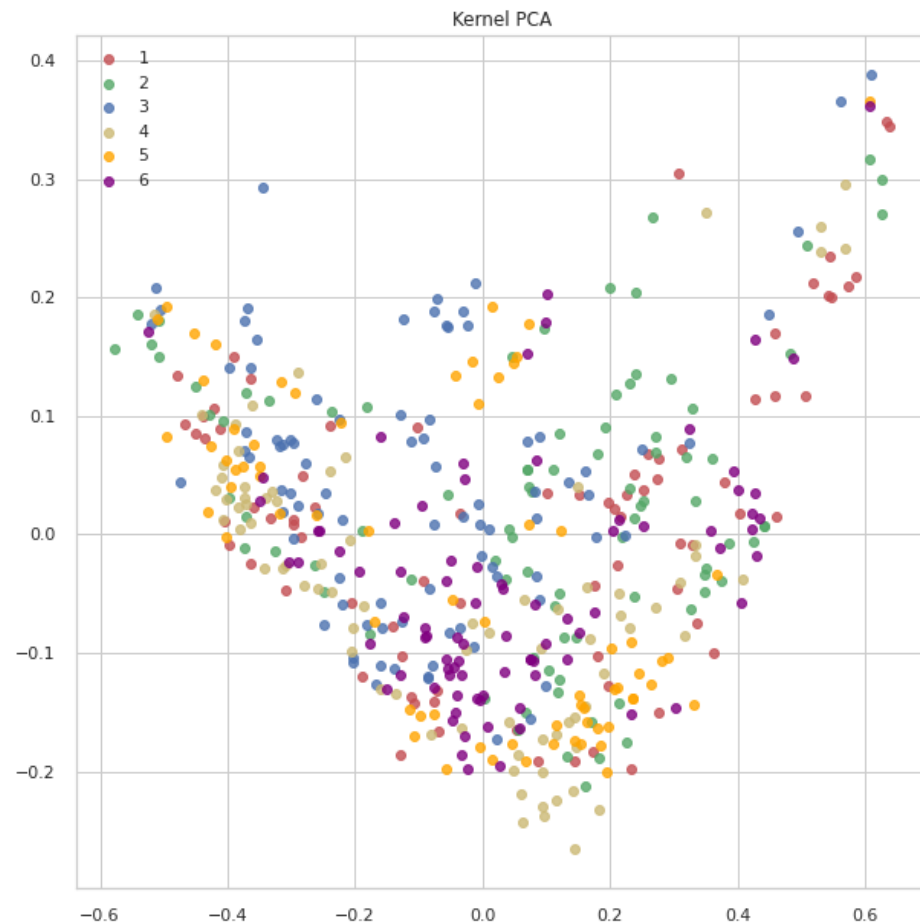
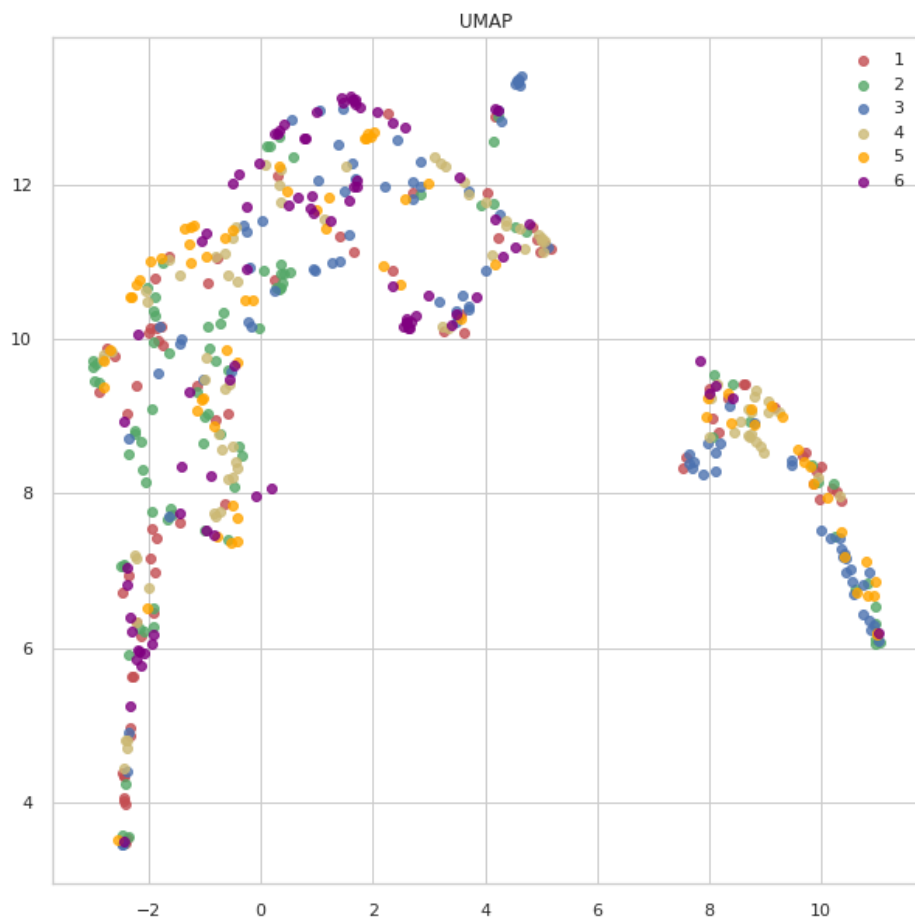
Visualizing outliers: unique structure and high/low number of edges/nodes



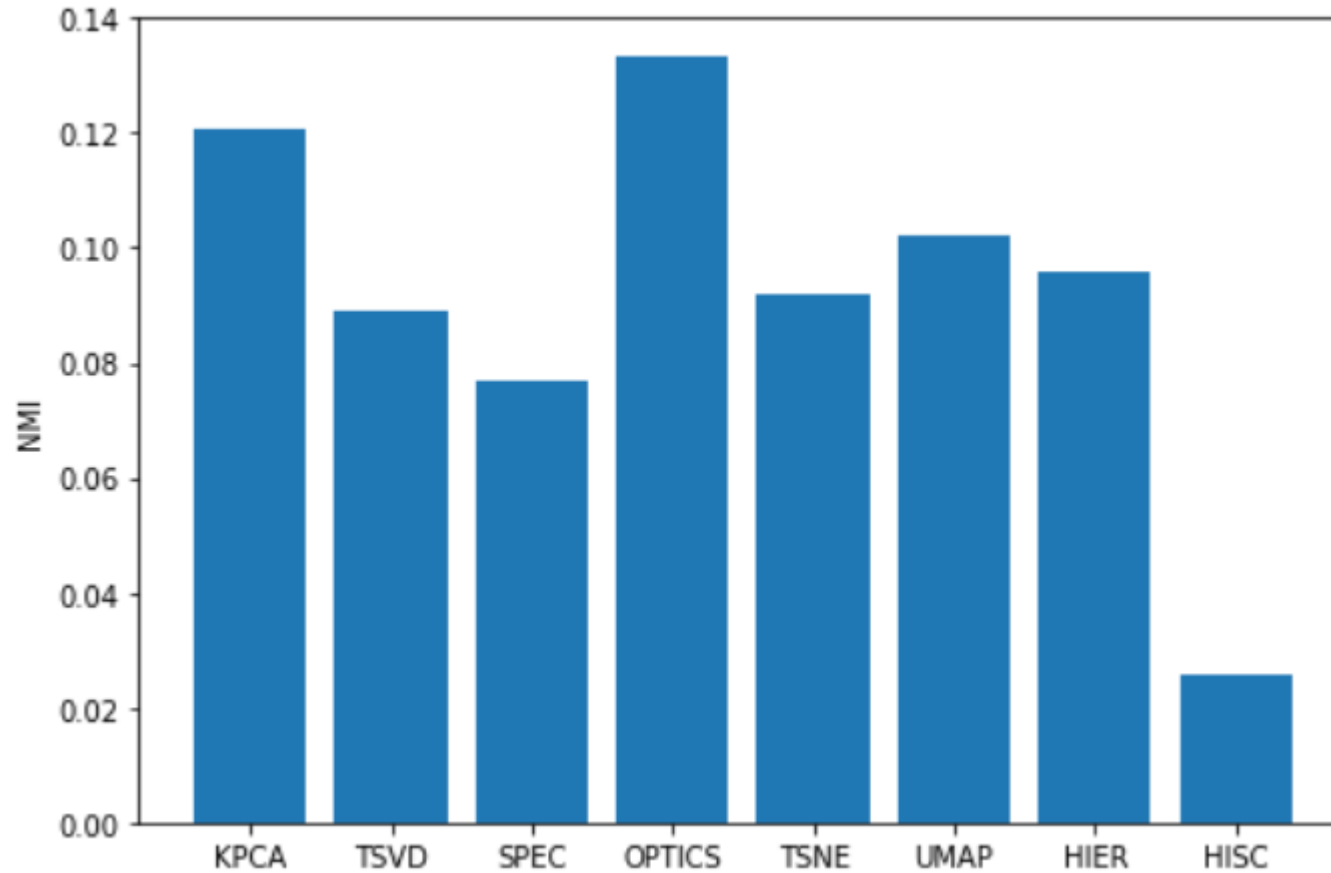
Preprocessing: removing outliers and graphs with extreme edge/node counts



Dimensionality Reduction on labeled dataset: WL4

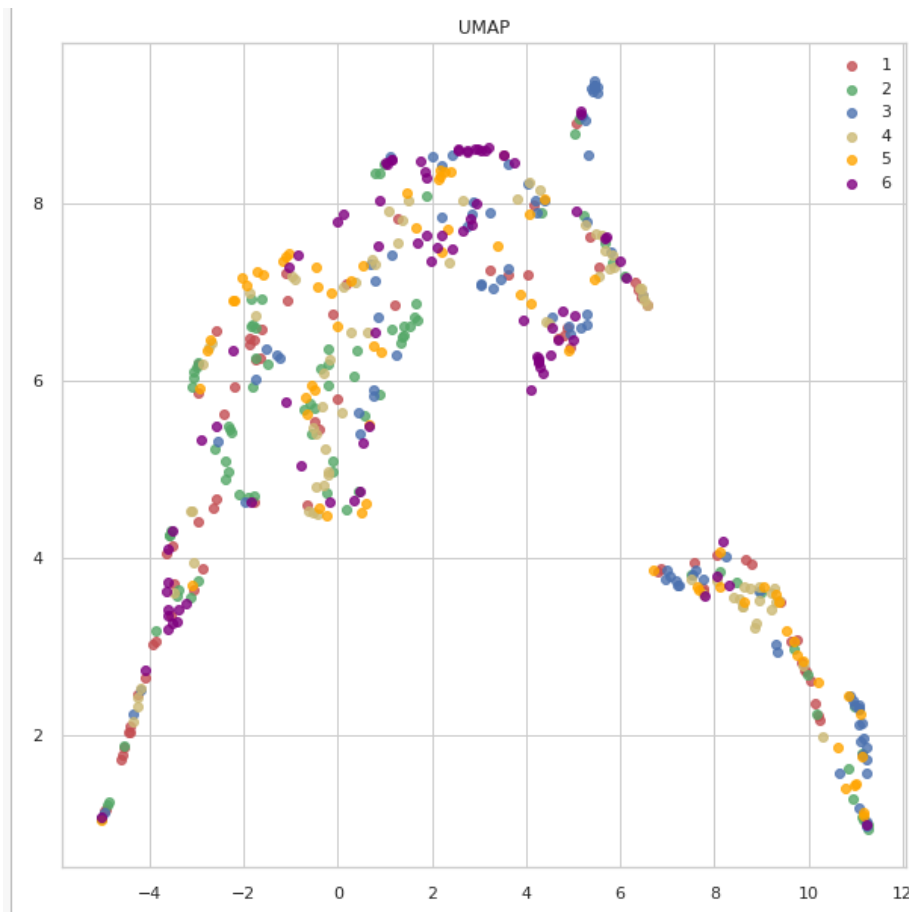


Cluster evaluation based on NMI

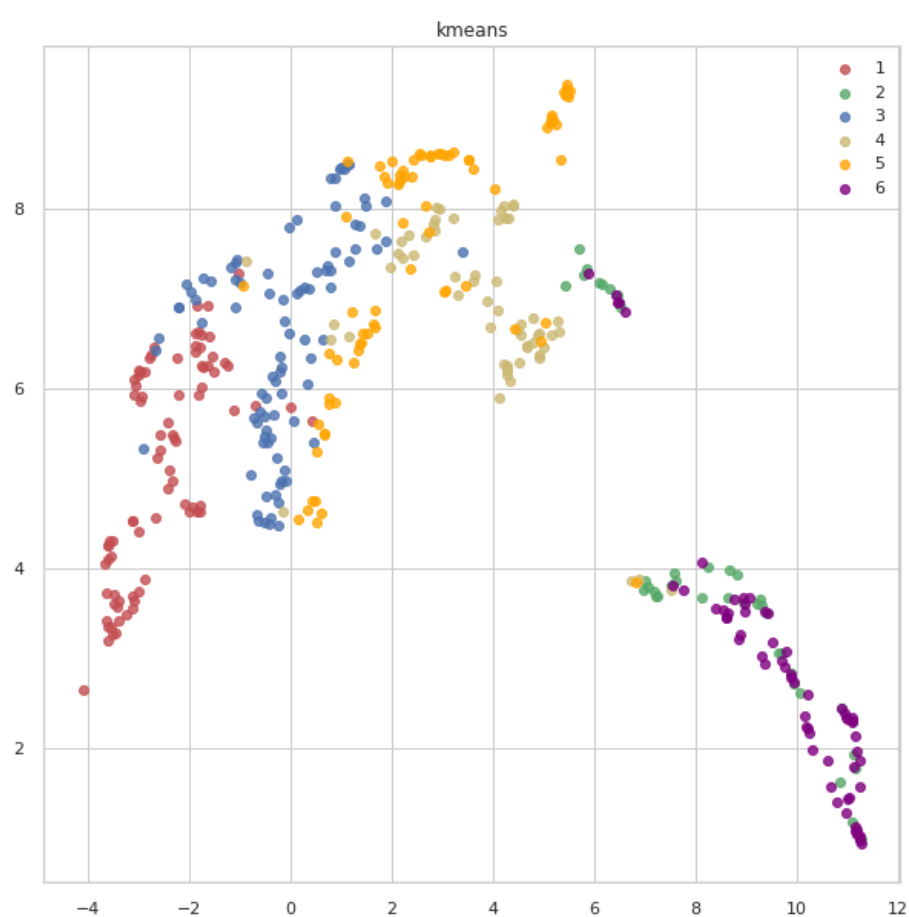


Clusters do not represent the true labels

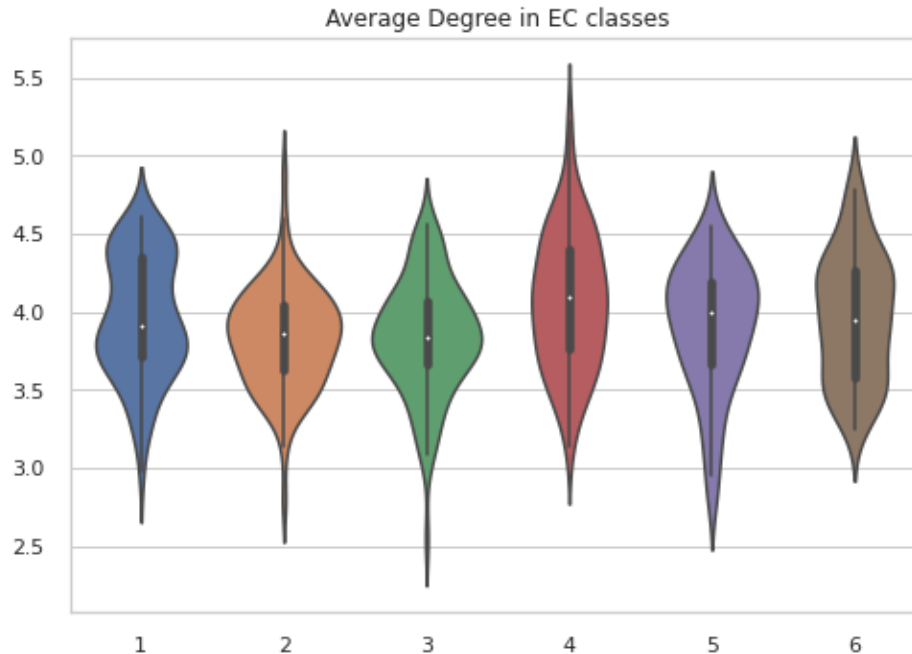
True labels



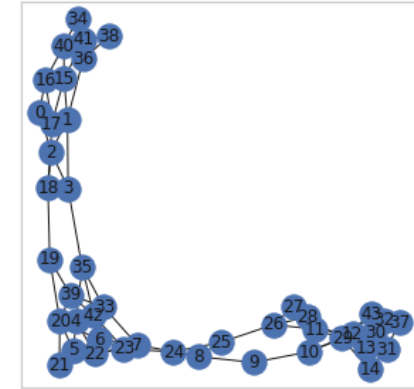
Kmeans Clusters



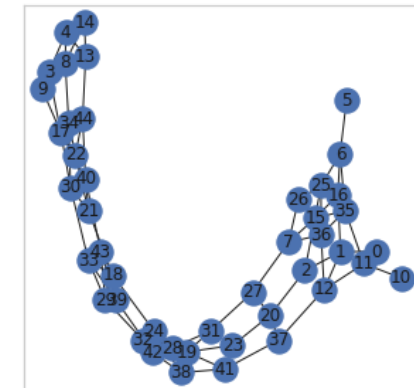
- Clusters have different degree distributions



Cluster 1, Indices [233 240]

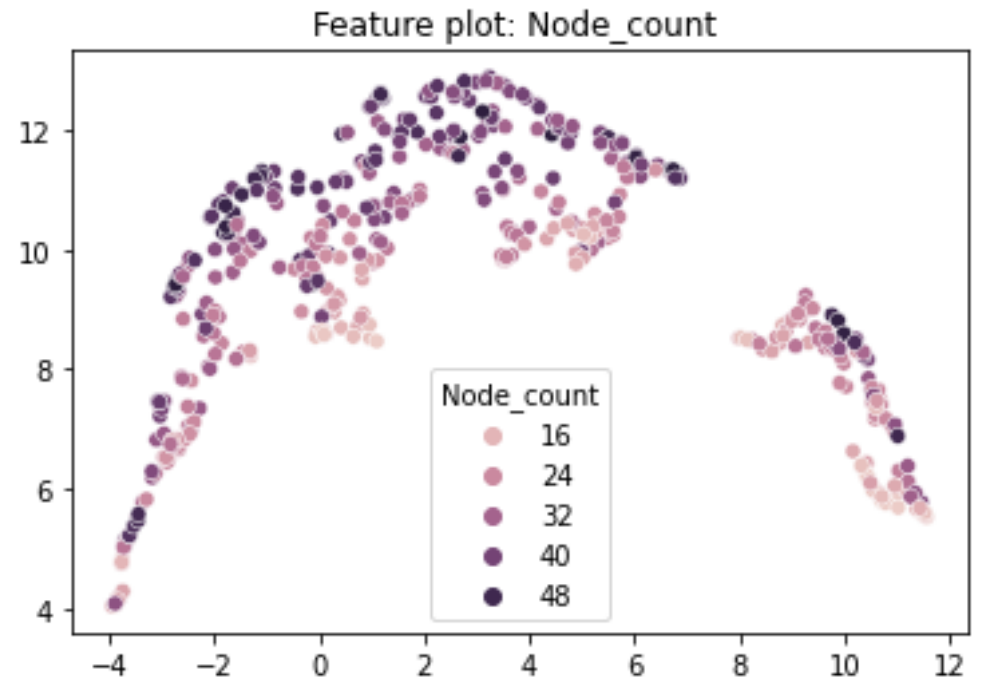
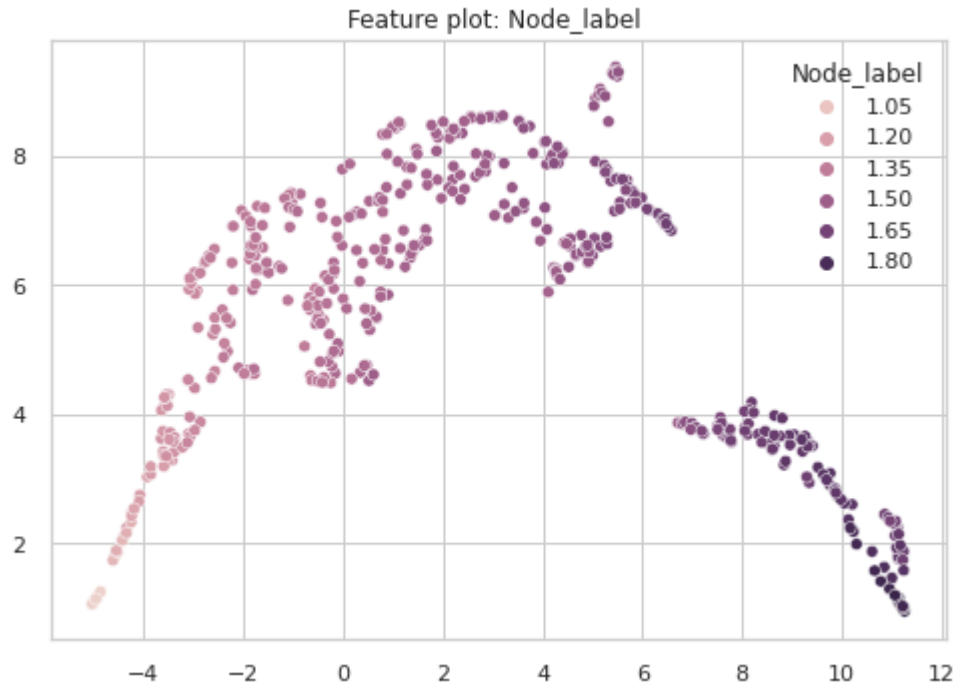


Cluster 2, Indices [485 232]



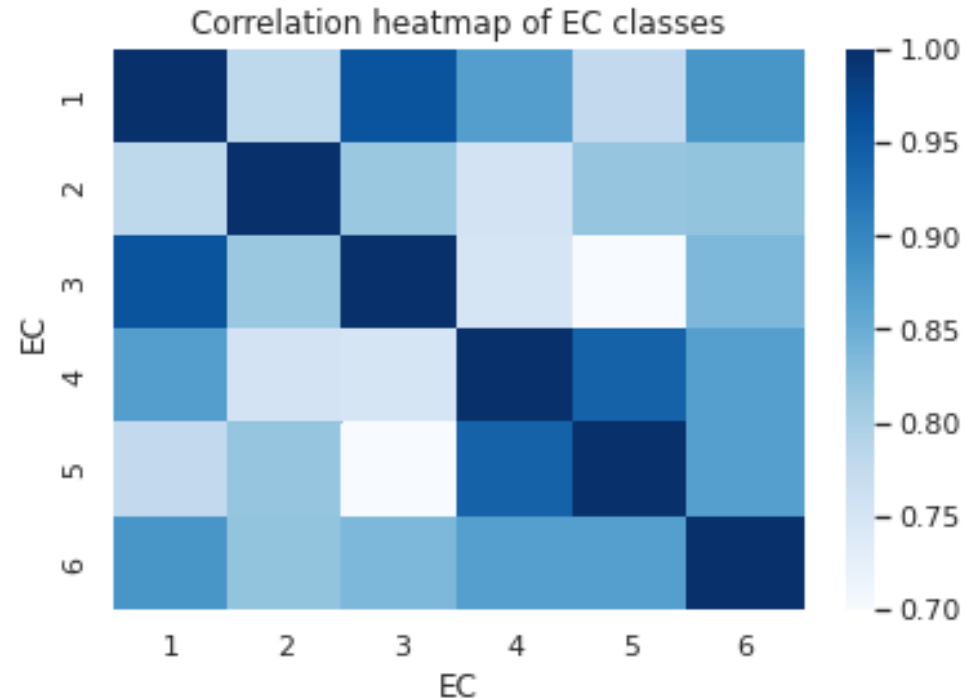
Group
8

DR describes node labels and node count

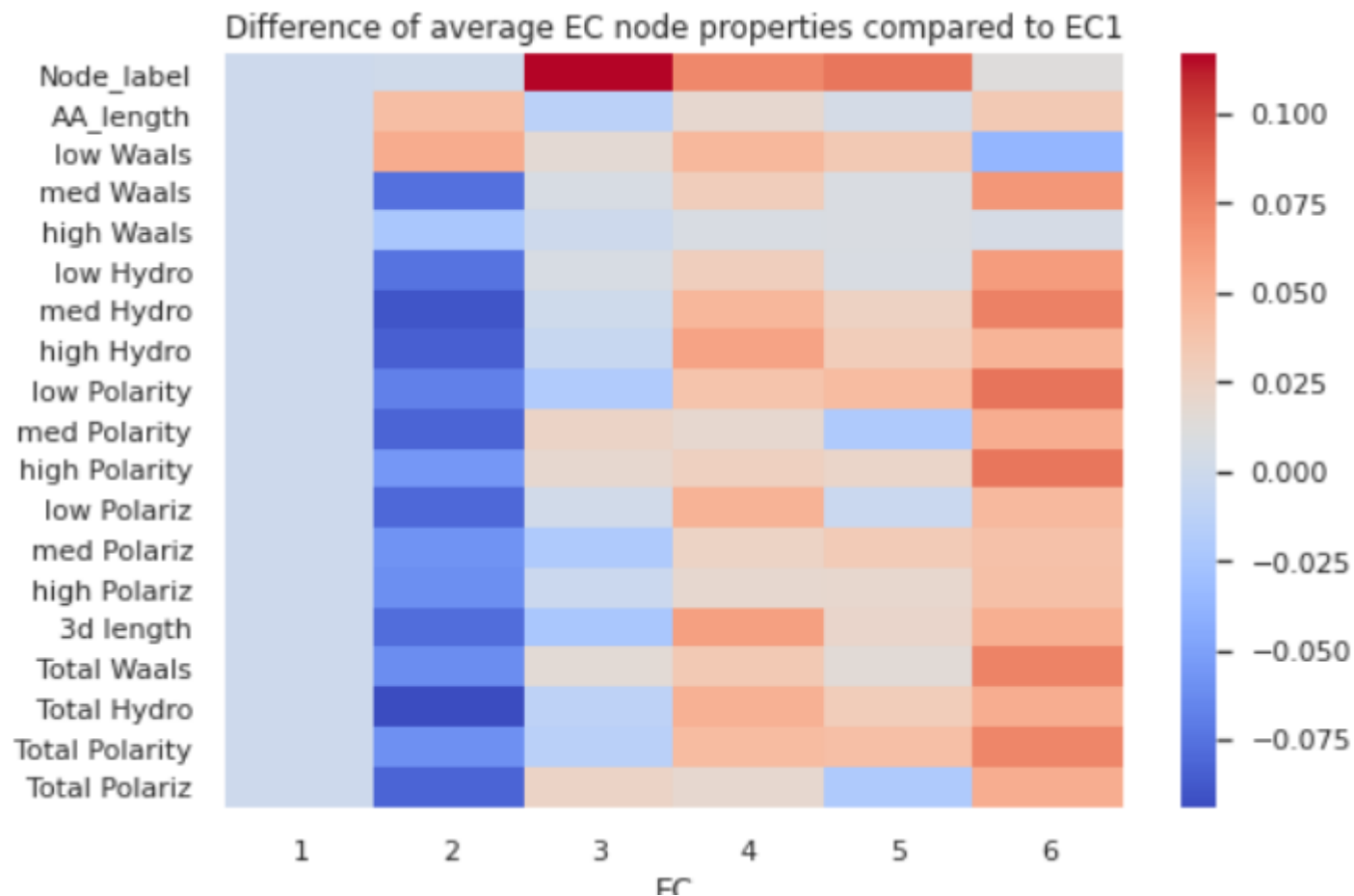


True label analysis based on node properties

- We notice higher correlation between EC 4,5
- Also EC 1,3 are highly correlated

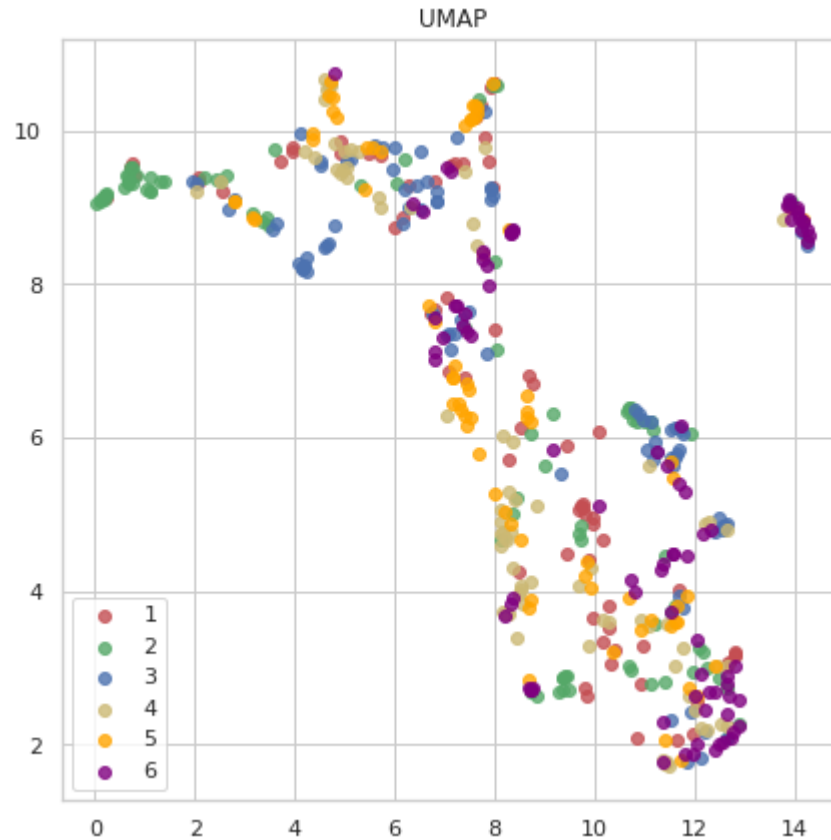


Difference between EC1 and other EC classes



Group
8

UMAP of only node attribute data (without labels)



Group
8

- SVM accuracy comparison of kernels

WL iterations 1-5 With preprocessing	0,44898
	0,56122
	0,57143
	0,62245
	0,61225
WL iterations 1-5 WithOUT preprocessing	0,45834
	0,575
	0,625
	0,625
	0,60834
Graphlet	0,32653
Shortest Paths	0,48979

Conclusion

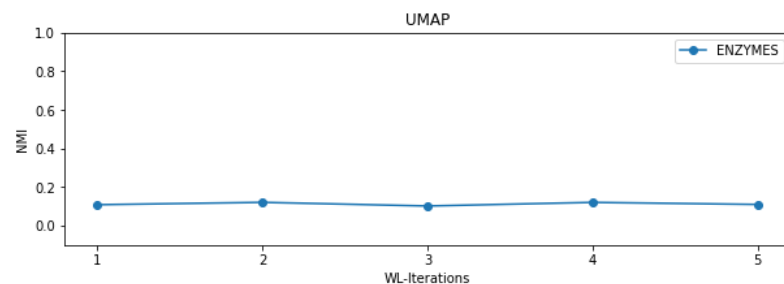
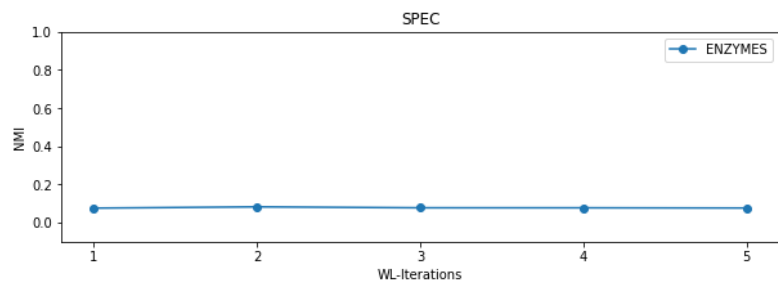
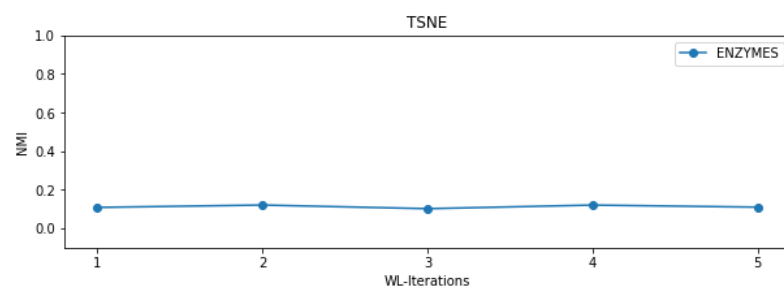
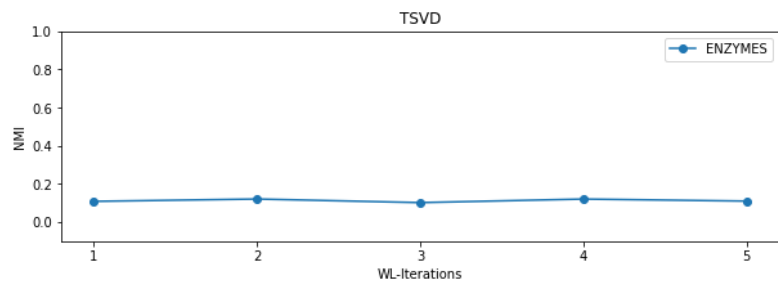
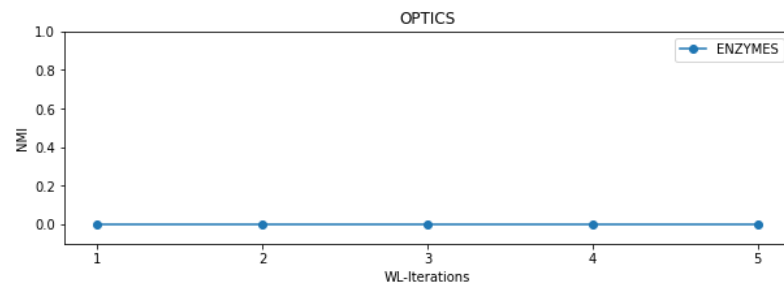
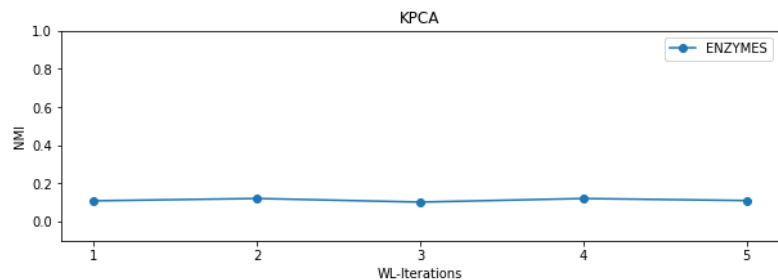
- UMAP vector representation is strongly driven by node labels (SSE elements), edge/node count
- Clusters do not represent the ground truth labels
- Clusters are formed based on different SSE composition, and structural differences e.g. degree distribution
- Node attributes play also an important role in EC classification which could be used for building better kernels

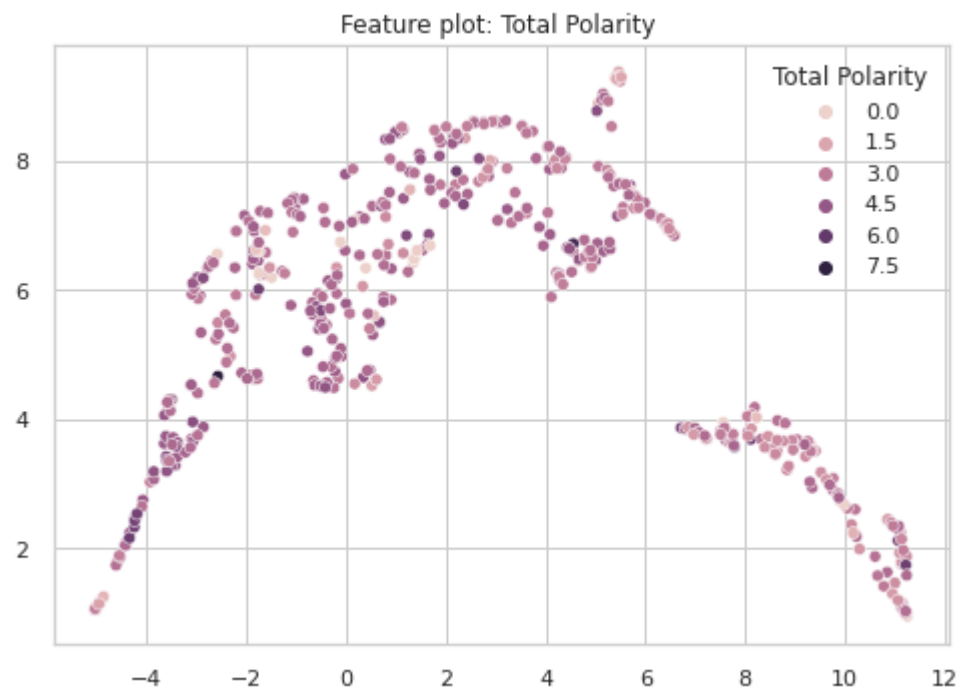
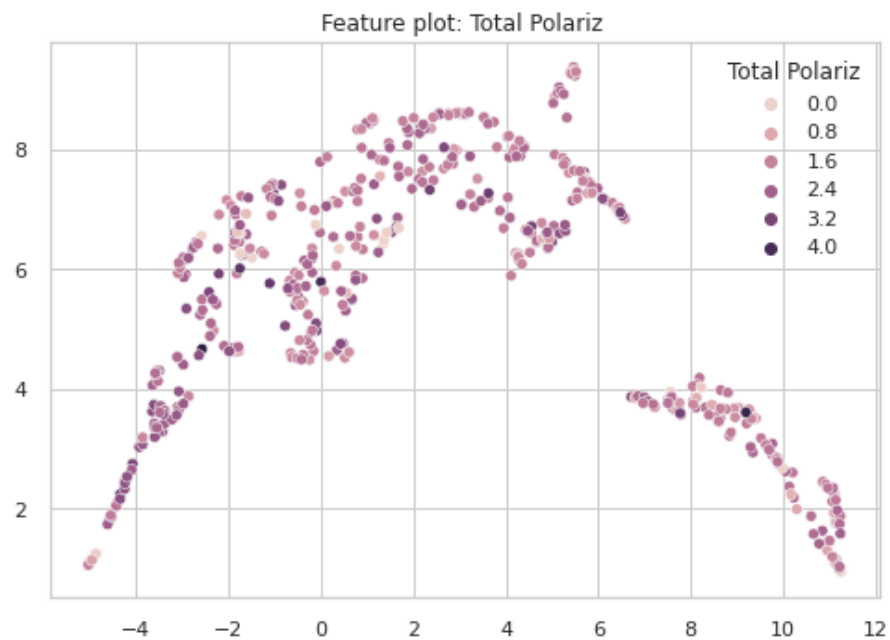
Thank you for your attention!

- NMI values for the 3 kernels (with removing the extreme values first)

	WL iter 3	Graphlet	Shortest Paths
Kmeans	0,1018	0,0828	0,1091
Hierarchical Ward	0,2079	0,1710	0,1866
Hierarchical Complete	0,1774	0,1576	0,1601
Hierarchical Average	0,0446	0,0887	0,0602
Hierarchical Single	0,0338	0,0551	0,0508
Subkmeans	0,0833	0,0717	0,0754

Overview of NMI performance for different representations





Group
0

Scatter plot of avg. nodes and avg. edges

