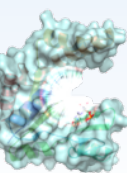


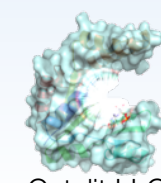
DataWknds.

CLUSTERING



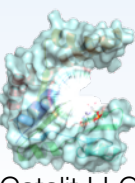
CLUSTERING

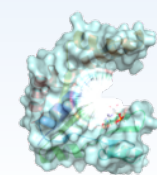
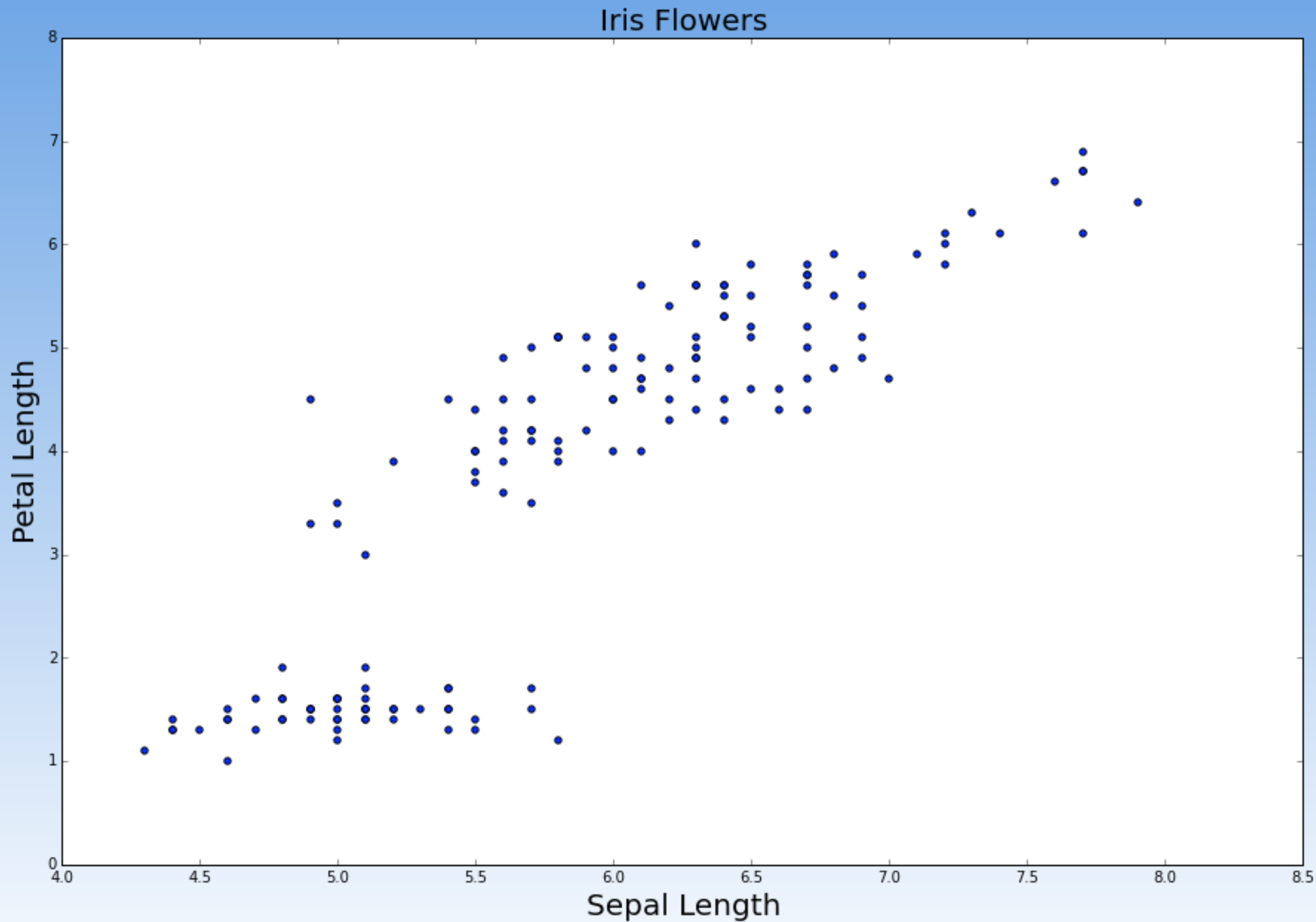
	CONTINUOUS	CATEGORICAL
SUPERVISED	?	?
UNSUPERVISED	?	?

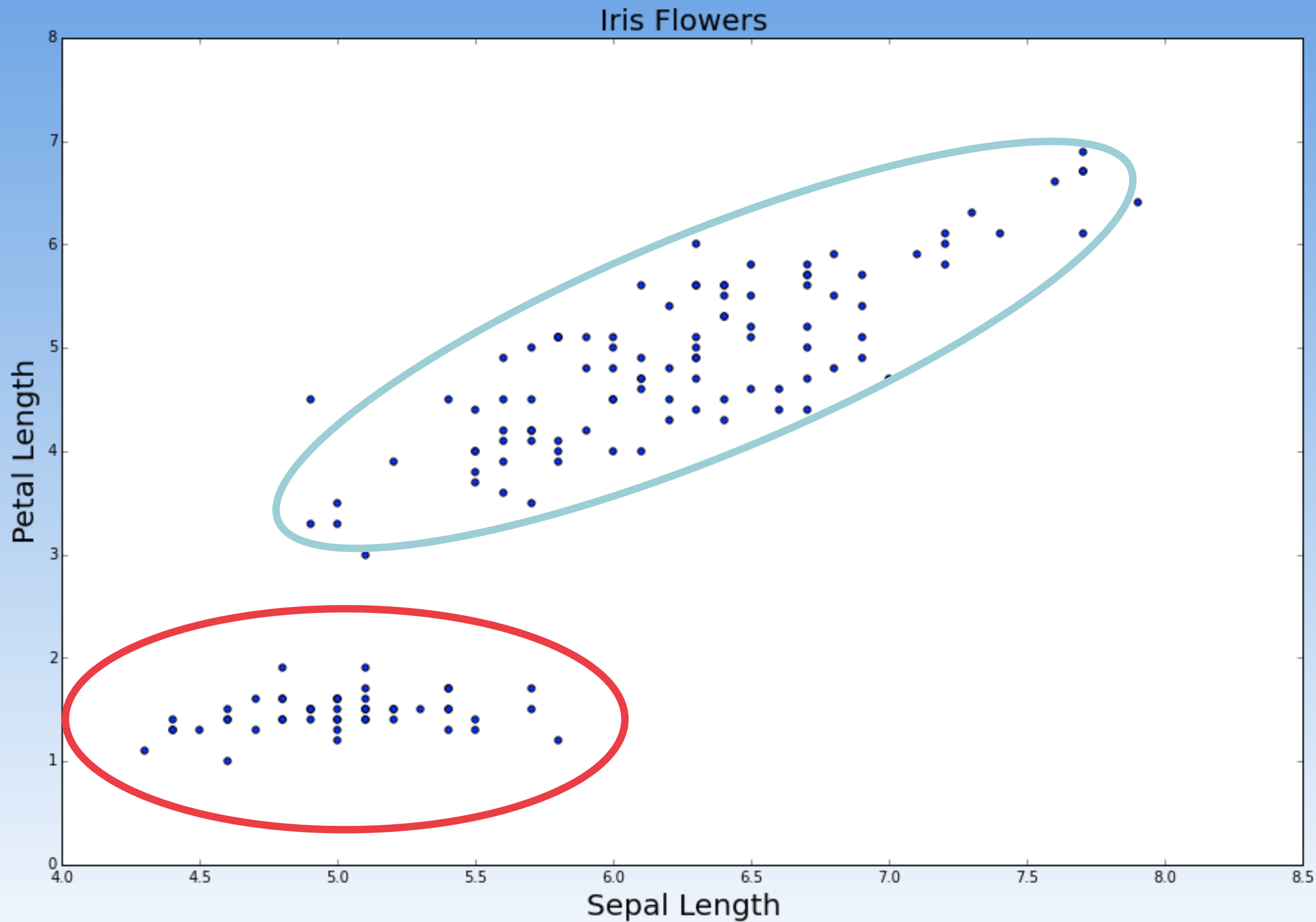


CLUSTERING

	CONTINUOUS	CATEGORICAL
SUPERVISED	REGRESSION	CLASSIFICATION
UNSUPERVISED	DIMENSION REDUCTION	CLUSTERING







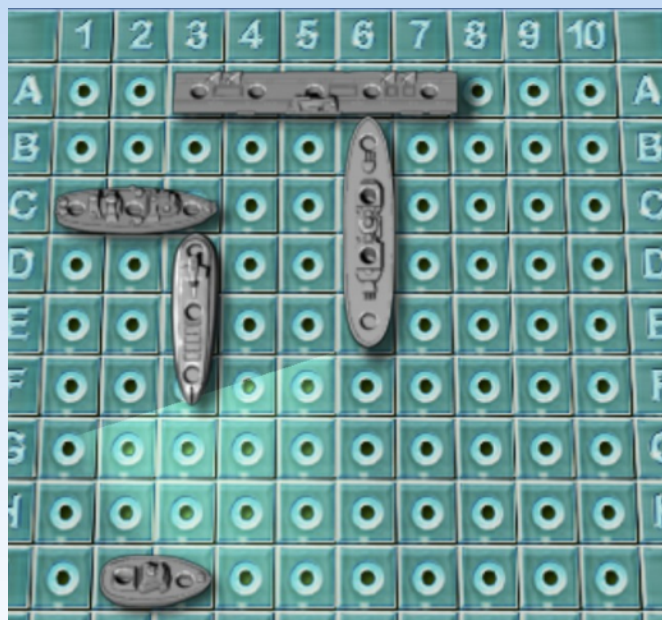
GOAL OF CLUSTERING

- GOAL: find data points that are similar to one another

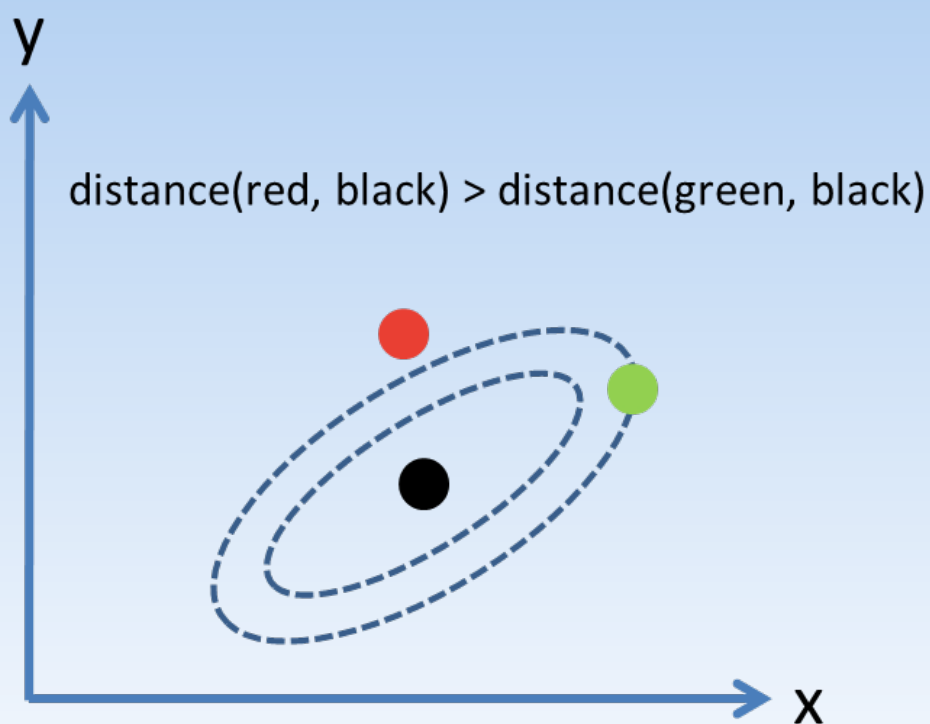
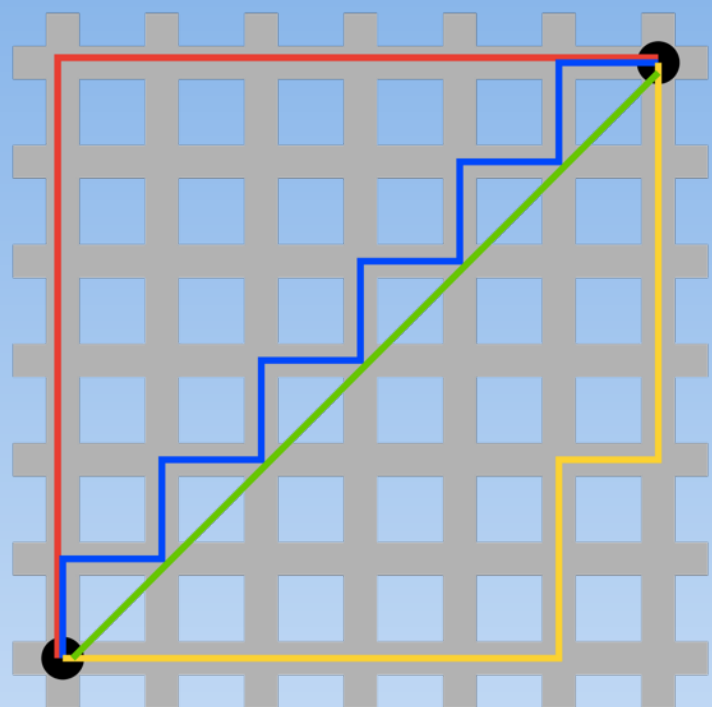


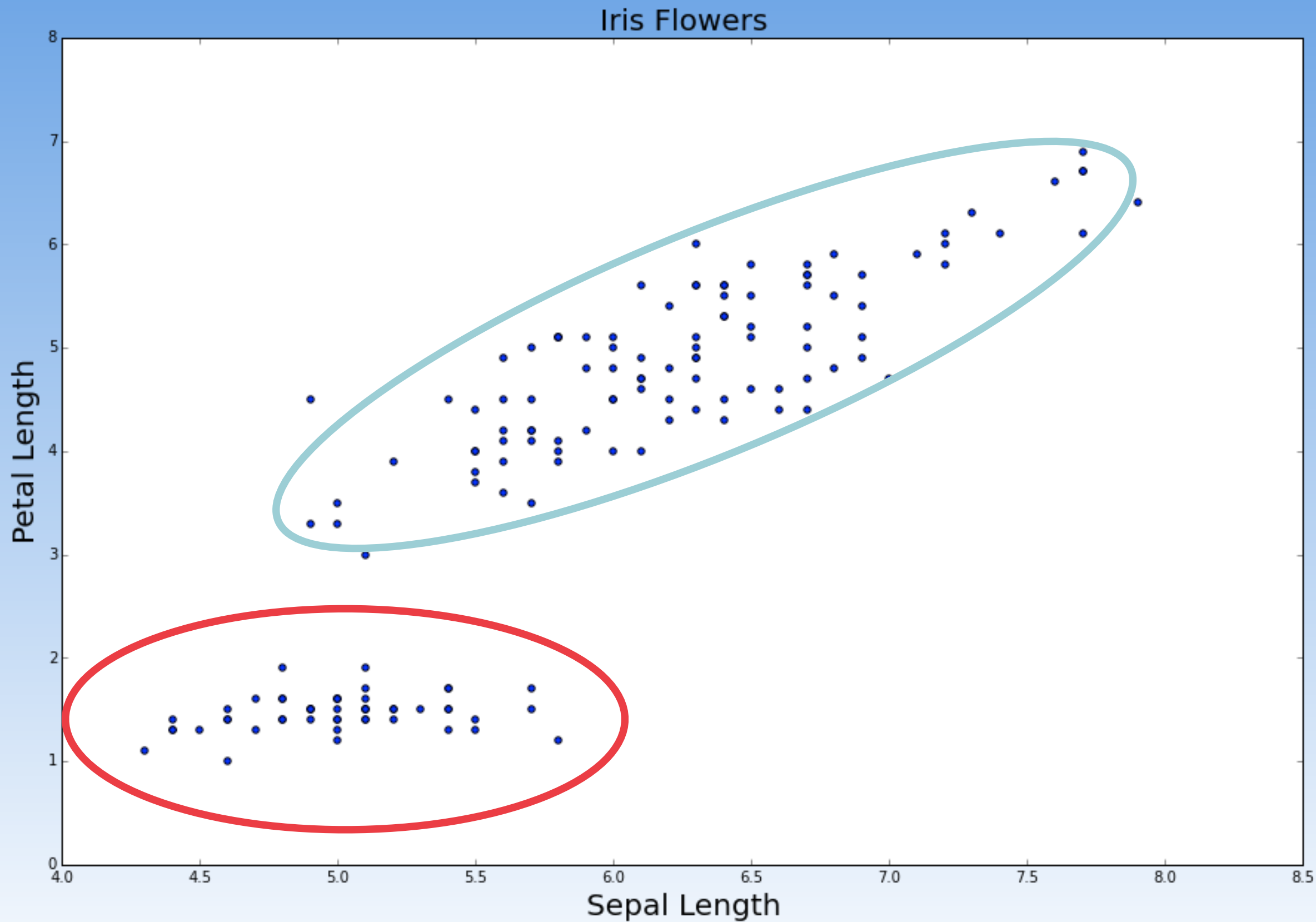
DISTANCE & SIMILARITY

	Age	Gender	Annual Salary	Months in residence	Months in job	Current Debt
Client 1	23	M	\$30,000	36	12	\$5,000
Client 2	30	F	\$45,000	12	12	\$1,000
Client 3	19	M	\$15,000	3	1	\$10,000



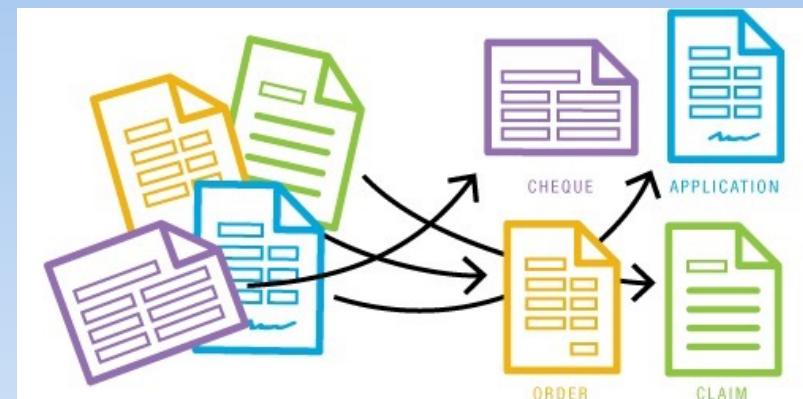
DISTANCE & SIMILARITY



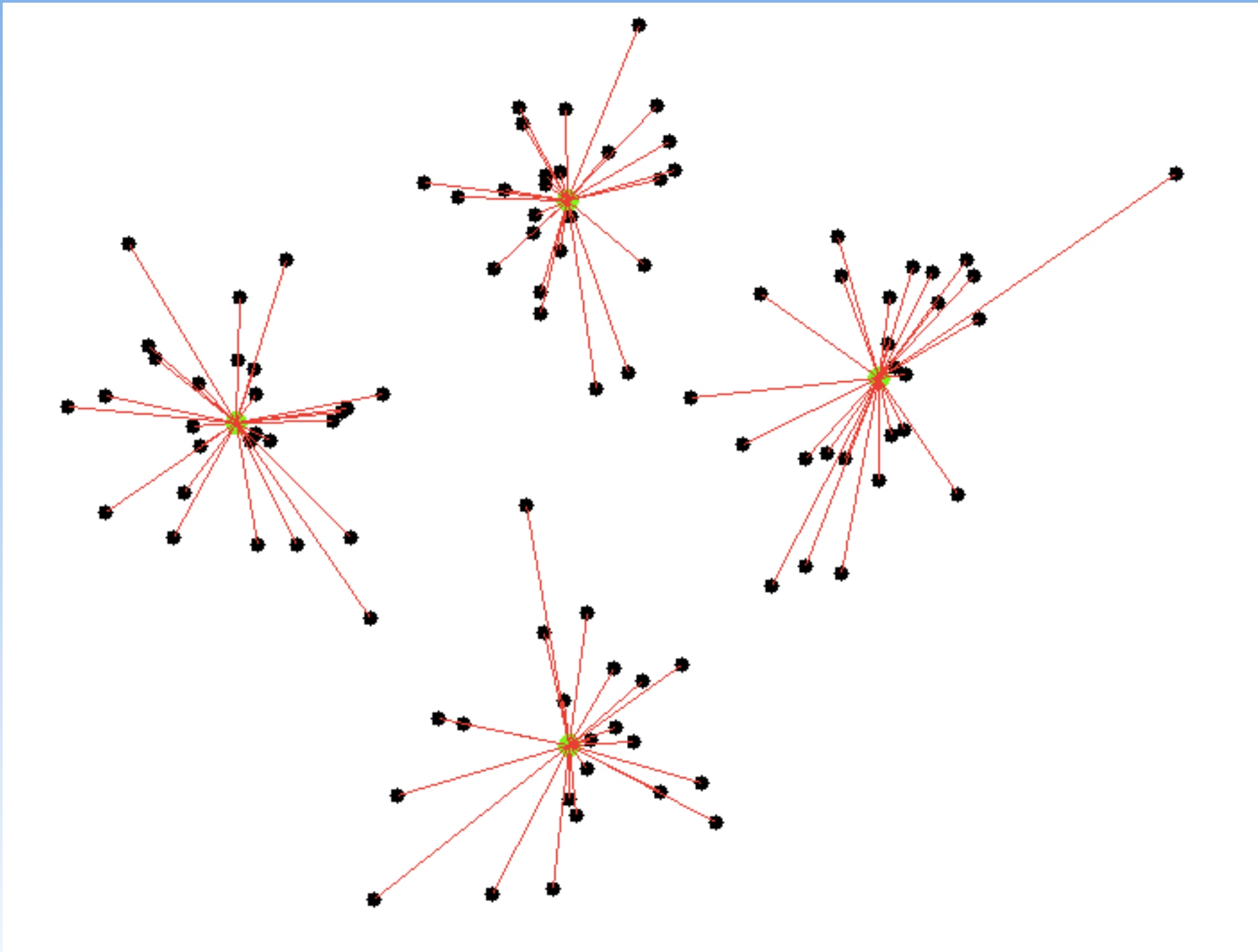


APPLICATIONS

- Malware detection
- Intrusion detection
- Document clustering
- Product placing
- Customer segmentation
- Gene discovery

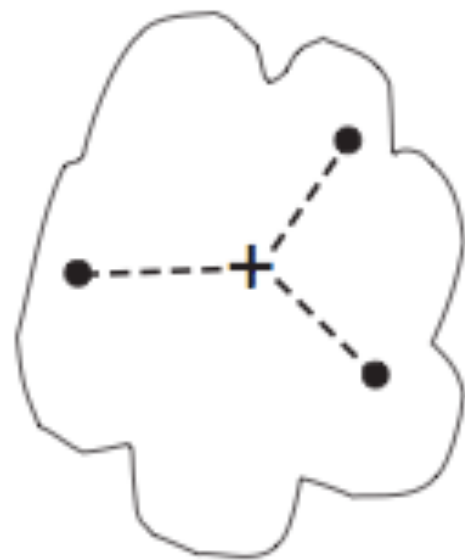


K-MEANS

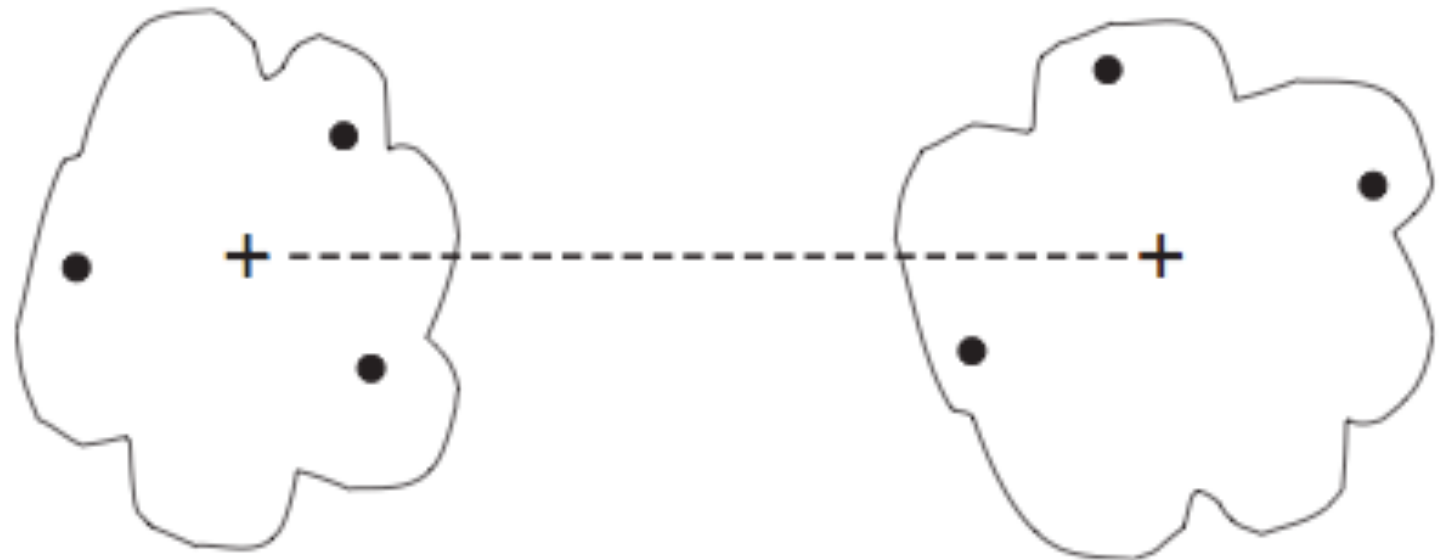


<http://shabal.in/visuals/kmeans/1.html>

CLUSTER VALIDATION



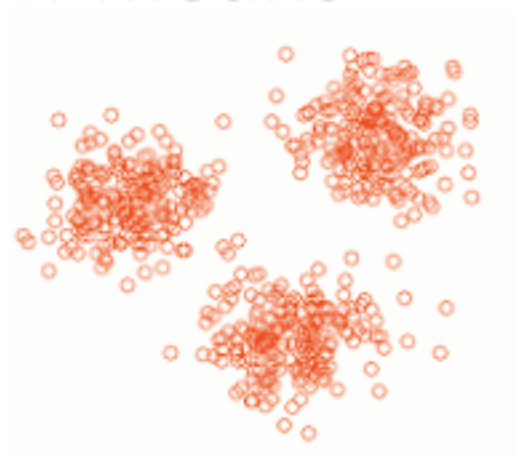
(a) Cohesion.



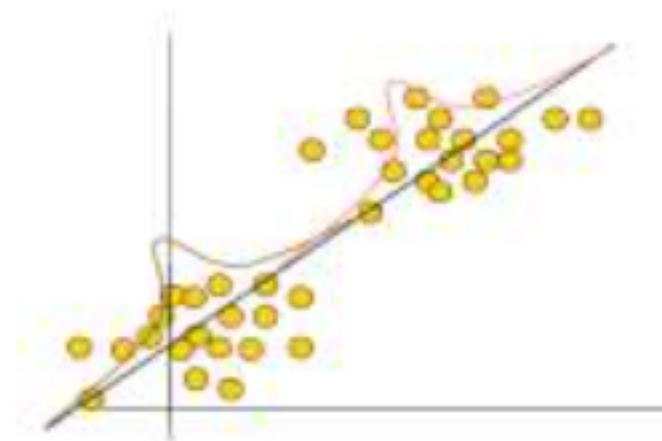
(b) Separation.

METHODS FOR CLUSTERING

- K-means



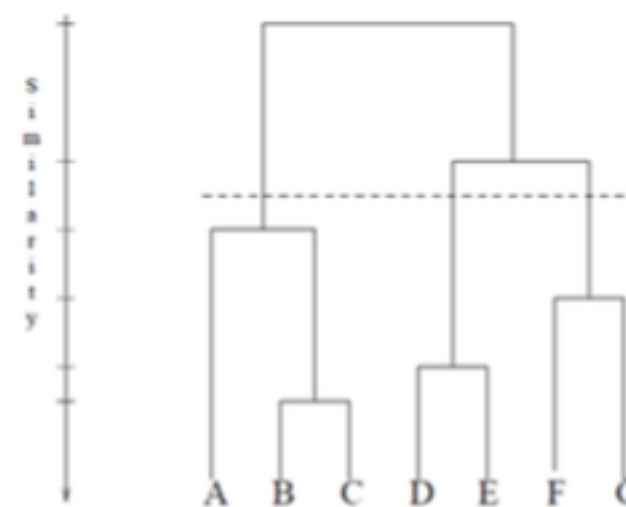
- > Expectation Maximization



- Density-based



- > Hierarchical



LAB: K-MEANS

