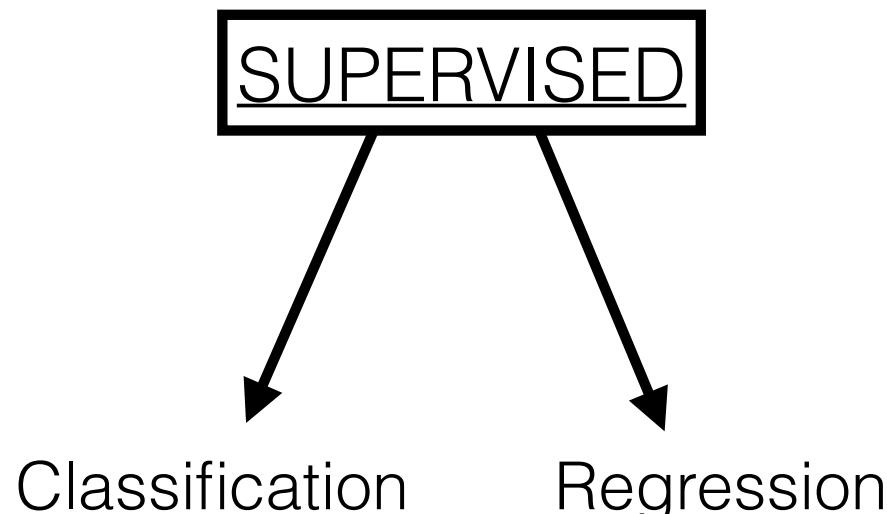


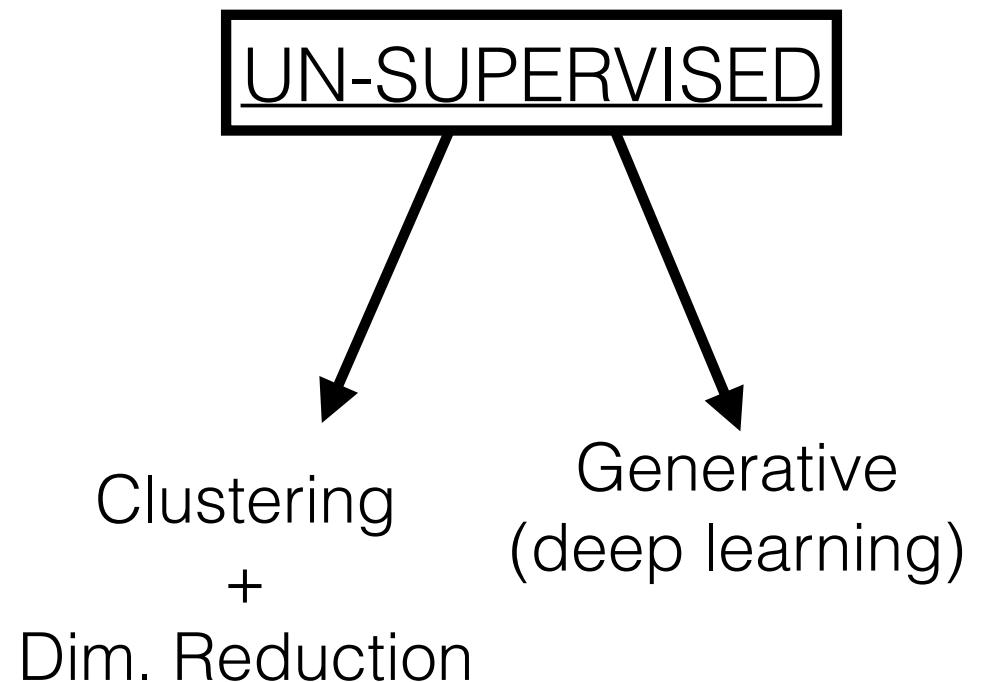
PART IV: A VERY BRIEF INTRODUCTION TO UNSUPERVISED LEARNING

WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for

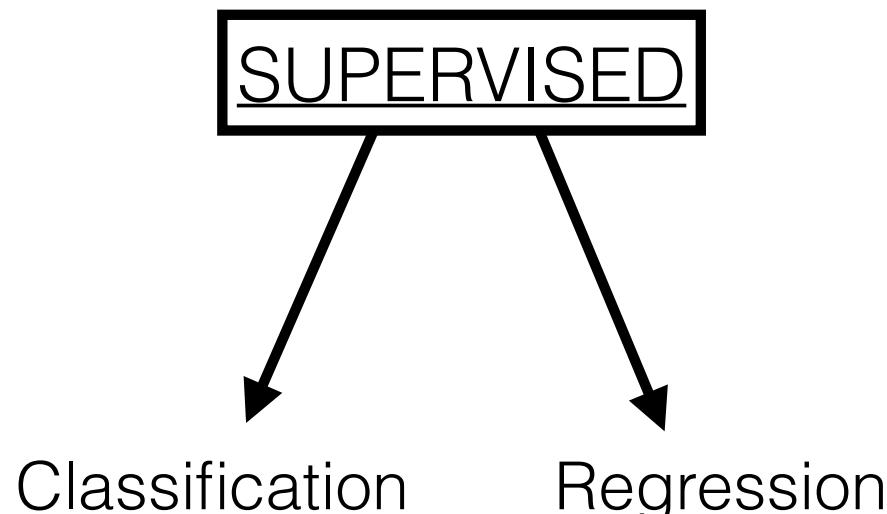


the machine is NOT told what to look for

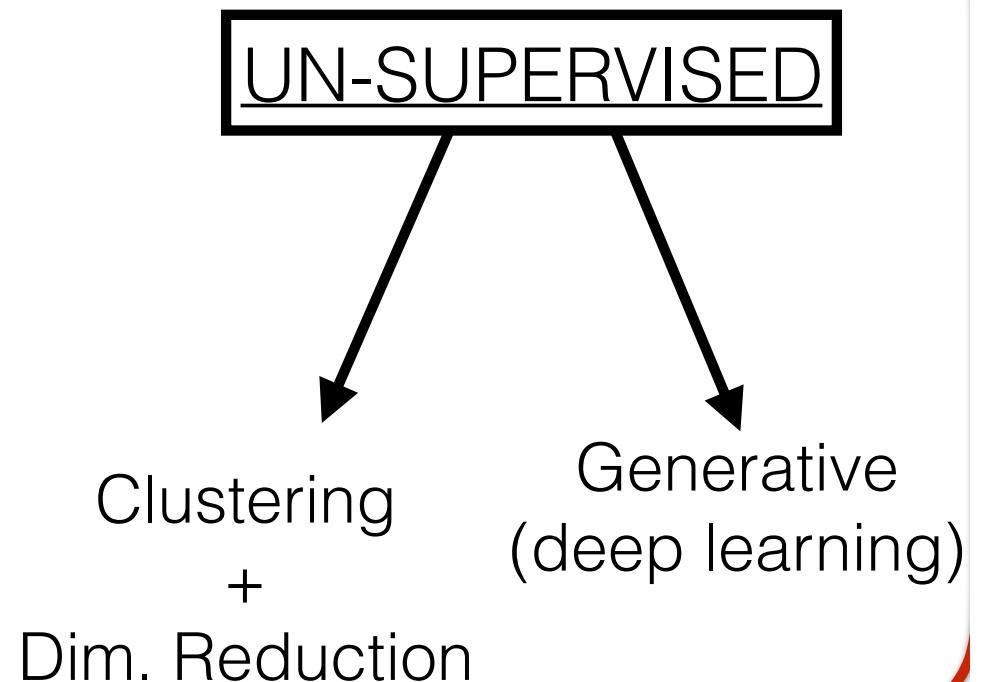


WHAT DOES MACHINE LEARNING DO?

the machine is told what to look for

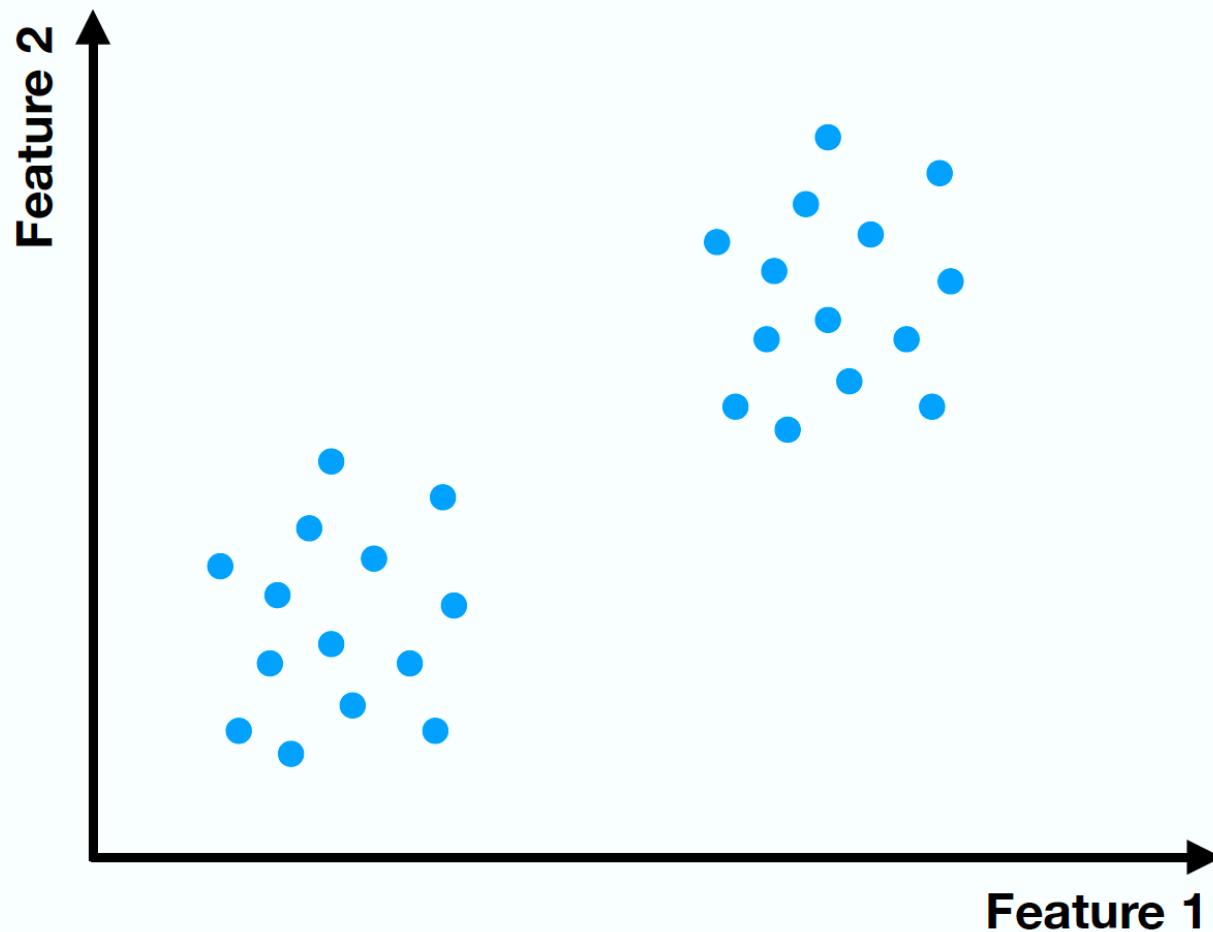


the machine is NOT told what to look for

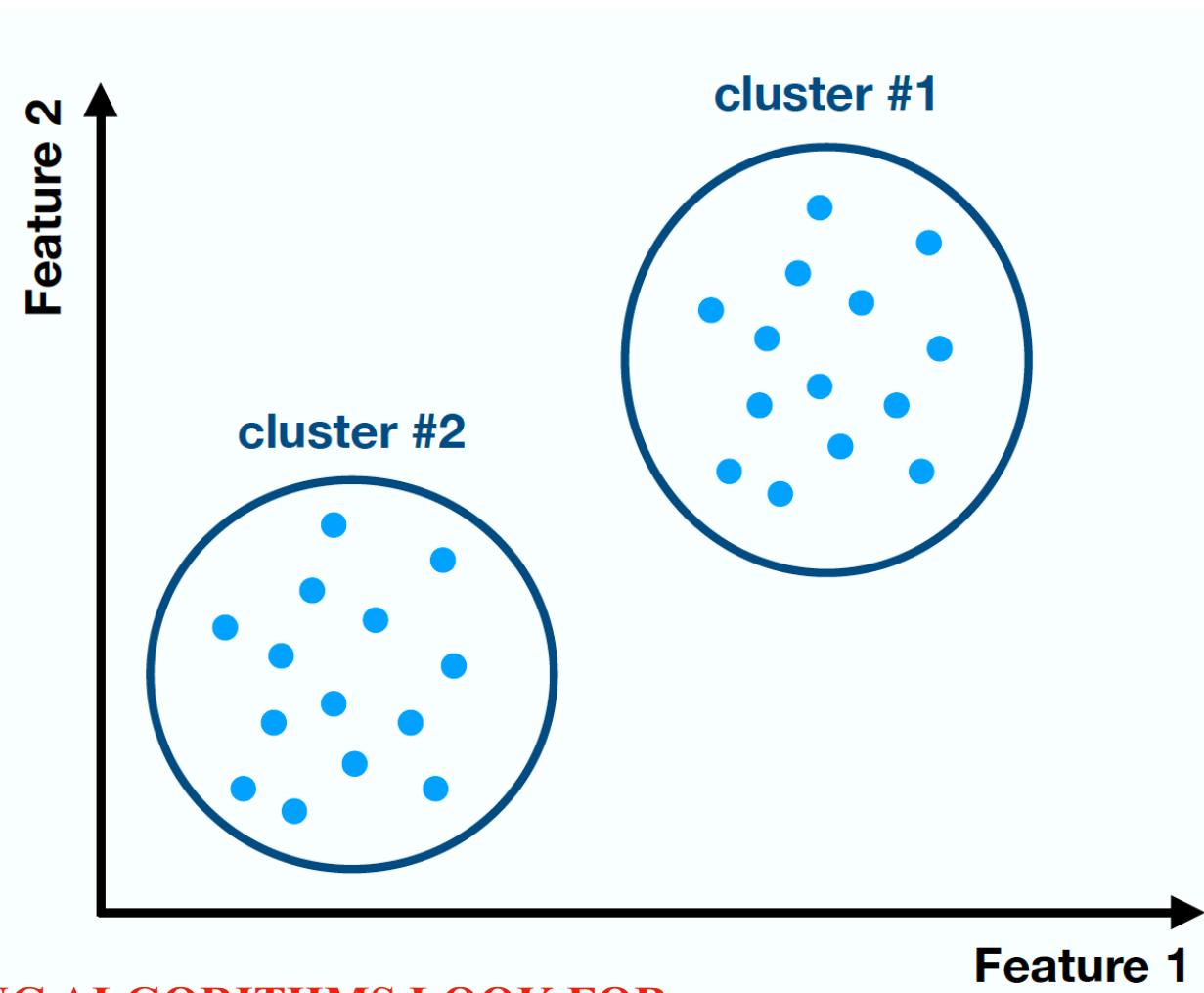


CLUSTERING ALGORITHMS

CLUSTERING ALGORITHMS



CLUSTERING ALGORITHMS



**CLUSTERING ALGORITHMS LOOK FOR
GROUPS OF DATA POINTS
IN THE DATA SPACE**

CLUSTERING ALGORITHMS

**WHAT ARE CLUSTERING
ALGORITHMS USEFUL FOR?**

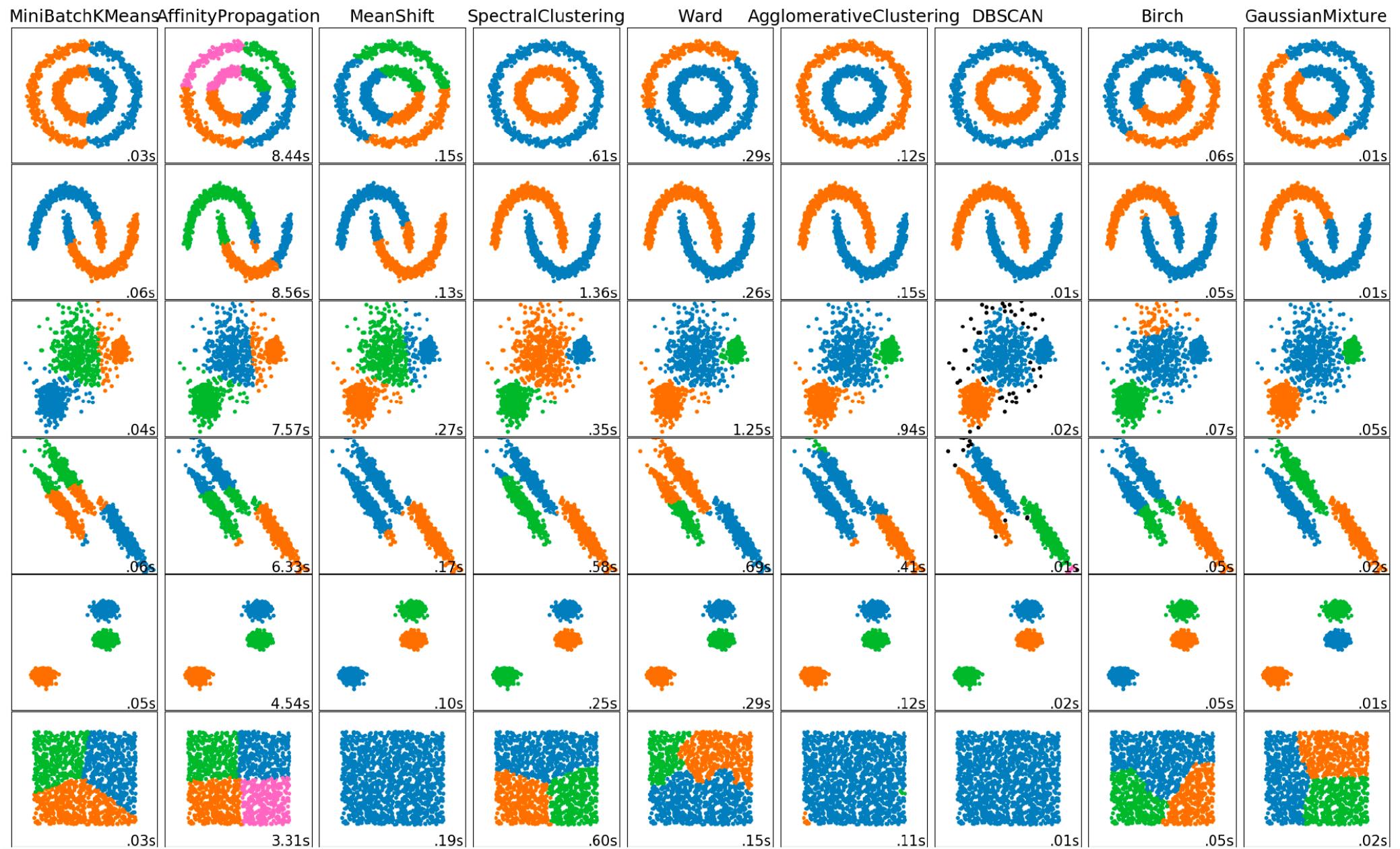
**1. CLASSIFICATION [OBJECTS WITH SIMILAR FEATURES ARE CLUSTERED
TOGETHER]**

2. VISUALIZATION OF LARGE DIMENSIONALITY SPACES

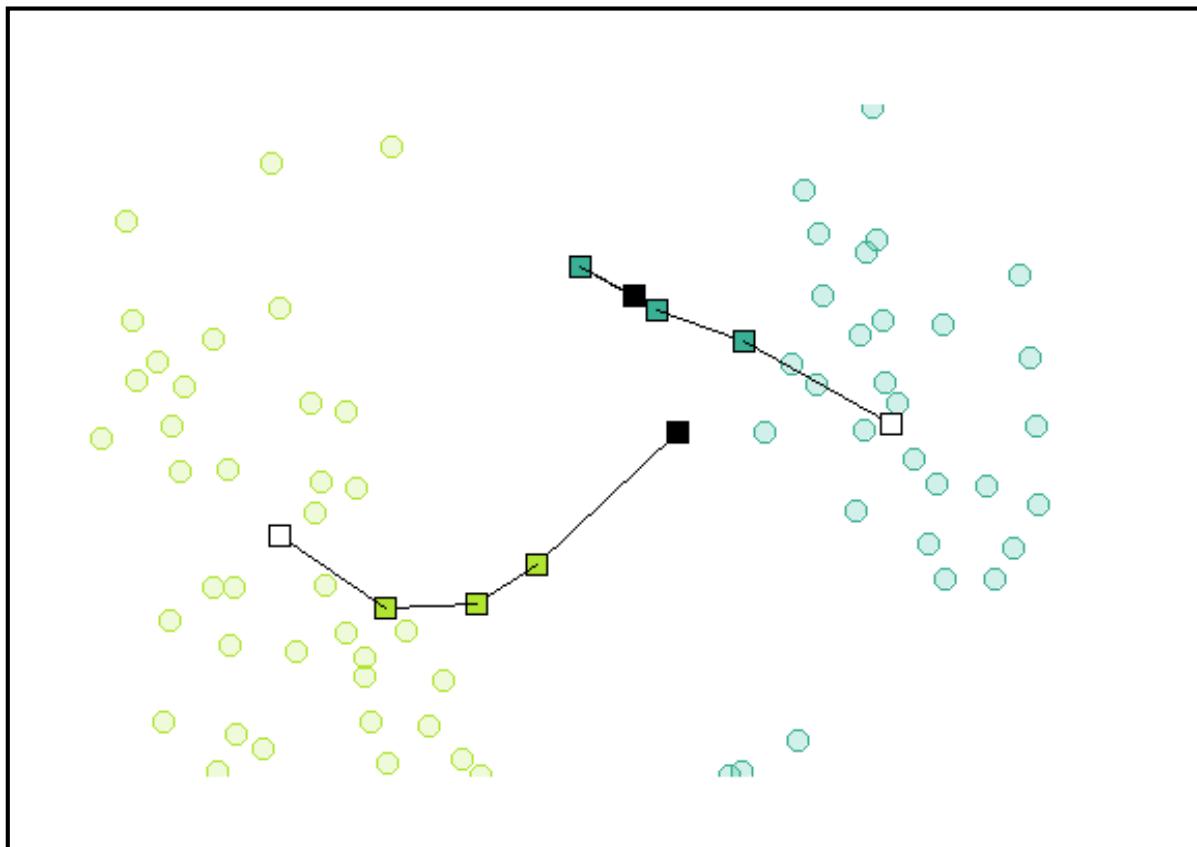
**CLUSTERING ALGORITHMS LOOK FOR
GROUPS OF DATA POINTS
IN THE DATA SPACE**

THERE ARE A LARGE NUMBER OF CLUSTERING ALGORITHMS. THEI BEHAVIOR DIFFERS DEPENDING ON THE DATA STRUCTURE.

[SCKIT LEARN]



EXAMPLE: K-MEANS



K-MEANS IS VERY SENSITIVE TO SEVERAL CRITICAL CHOICES

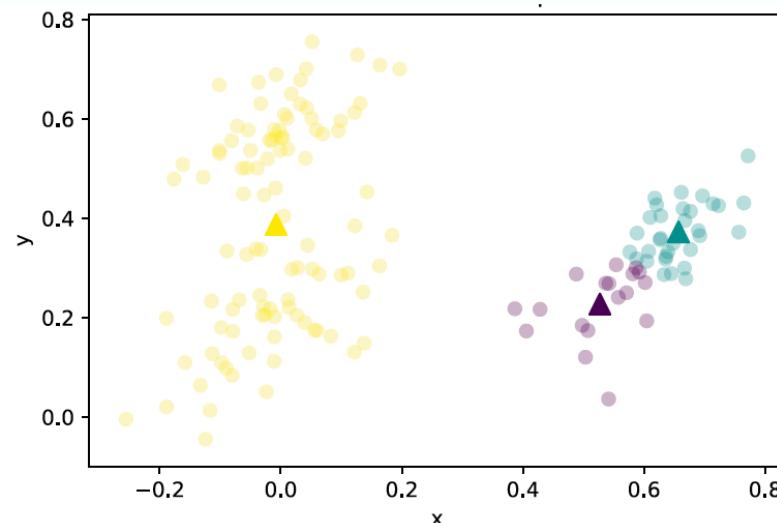
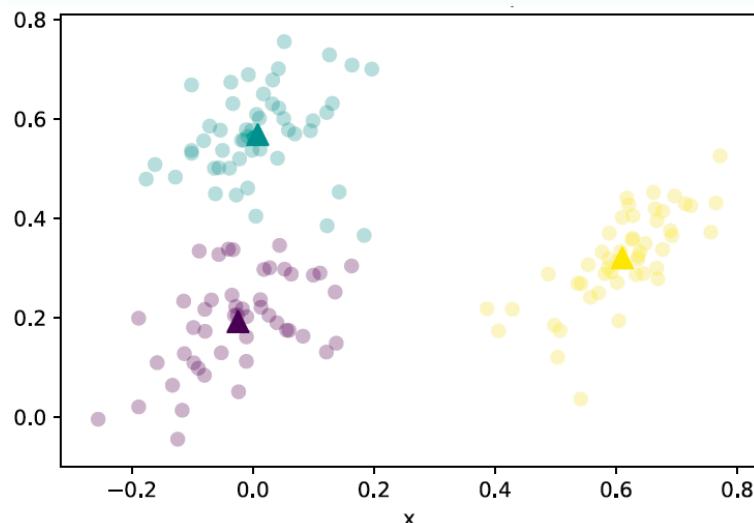
$$f(\vec{x}, \{a_1, a_2, \dots\}) = \vec{y}$$

Initial conditions:

1. Distance metric
2. Initial centroids

$$L = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

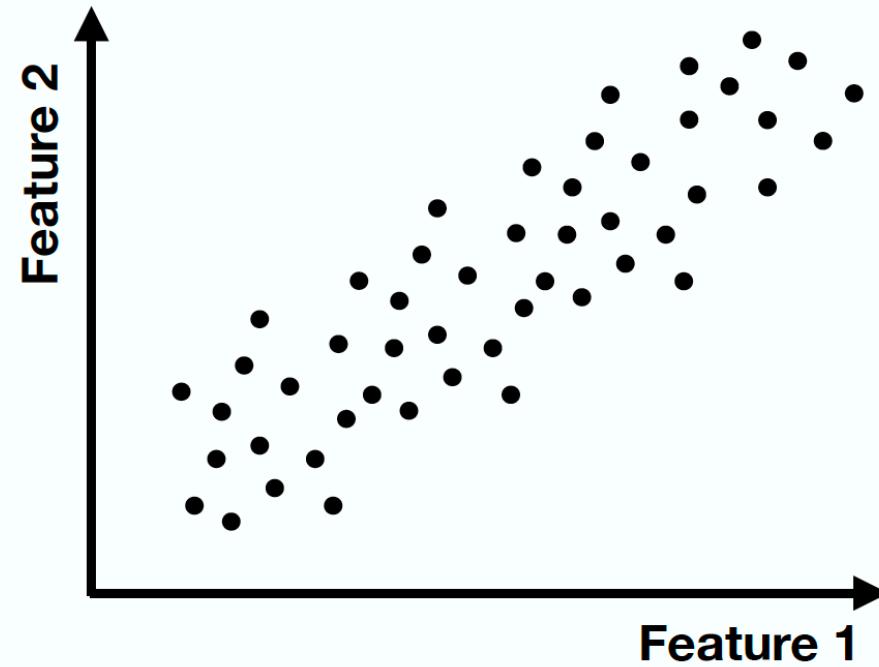
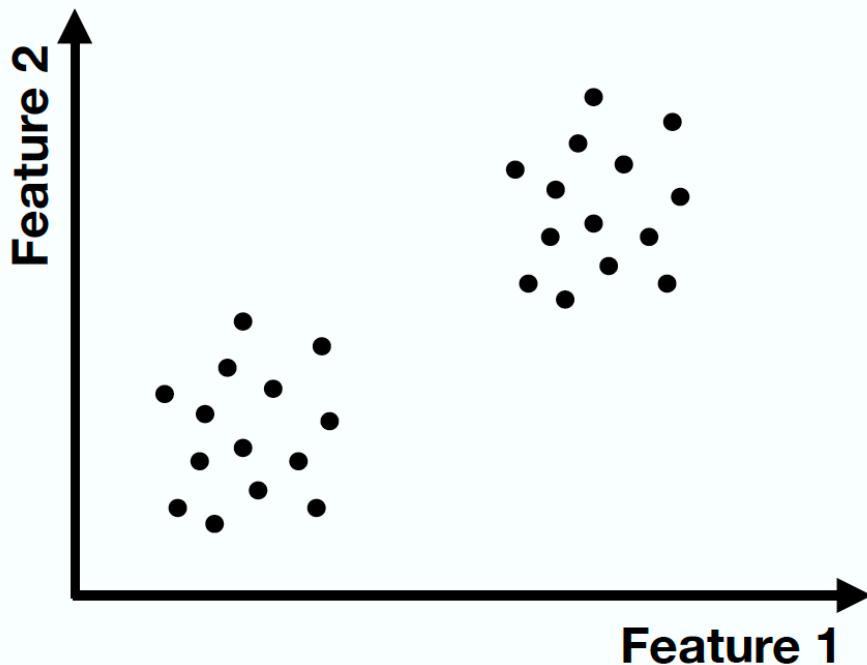
$k=3$, and two different random placements of centroids



K-MEANS IS VERY SENSITIVE TO SEVERAL CRITICAL CHOICES

$$f(\vec{x}, \{a_1, a_2, \dots\}) = \vec{y}$$

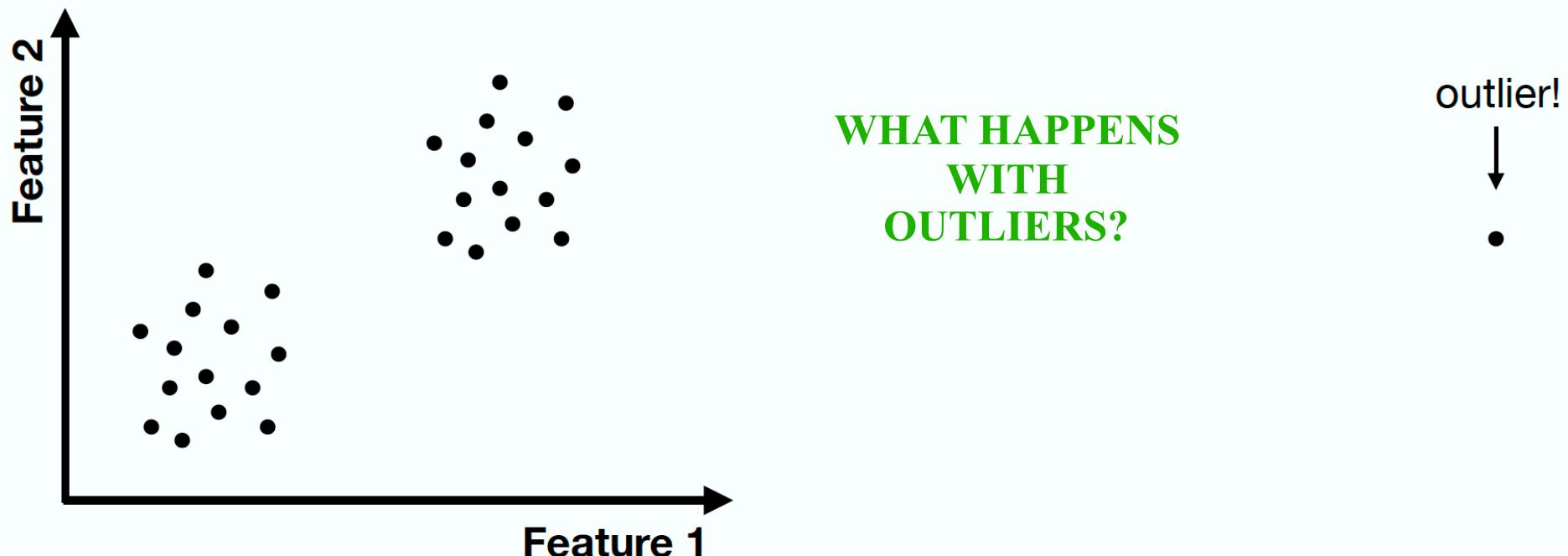
The behavior will strongly depend on the structure of the input features



K-MEANS IS VERY SENSITIVE TO SEVERAL CRITICAL CHOICES

$$f(\vec{x}, \{a_1, a_2, \dots\}) = \vec{y}$$

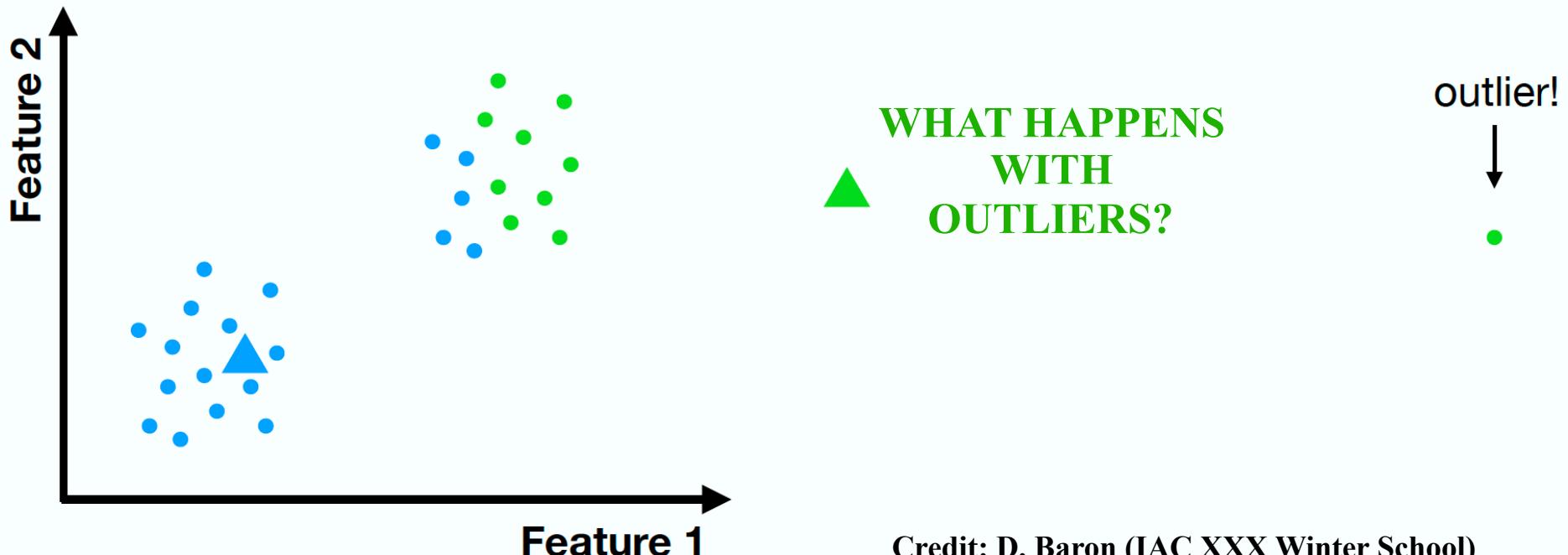
The behavior will strongly depend on the structure of the input features



K-MEANS IS VERY SENSITIVE TO SEVERAL CRITICAL CHOICES

$$f(\vec{x}, \{a_1, a_2, \dots\}) = \vec{y}$$

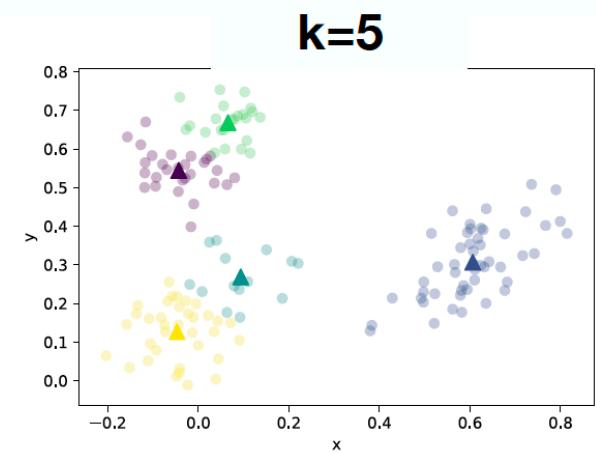
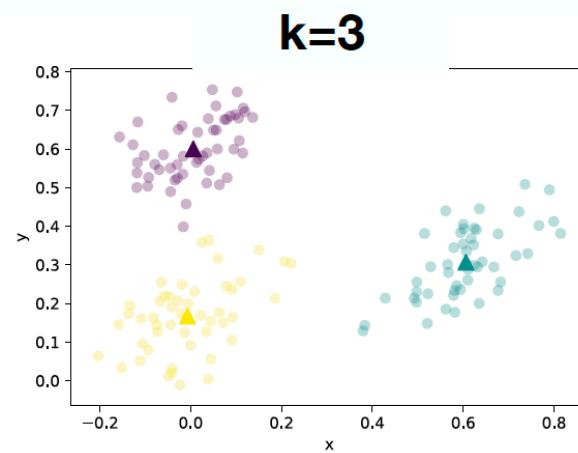
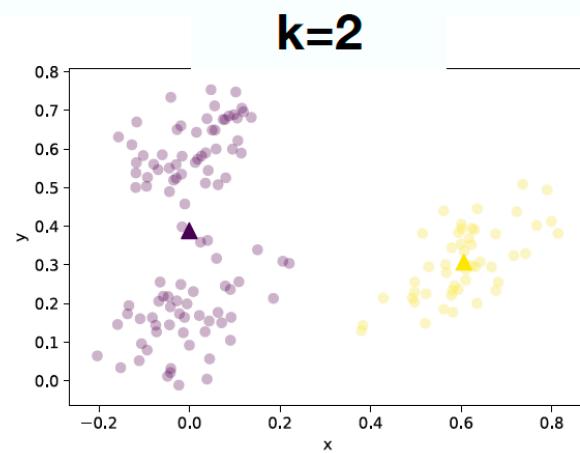
The behavior will strongly depend on the structure of the input features



K-MEANS IS VERY SENSITIVE TO SEVERAL CRITICAL CHOICES

$$f(\vec{x}, \{a_1, a_2, \dots\}) = \vec{y}$$

Hyper-Parameters: The number of clusters needs to be fixed

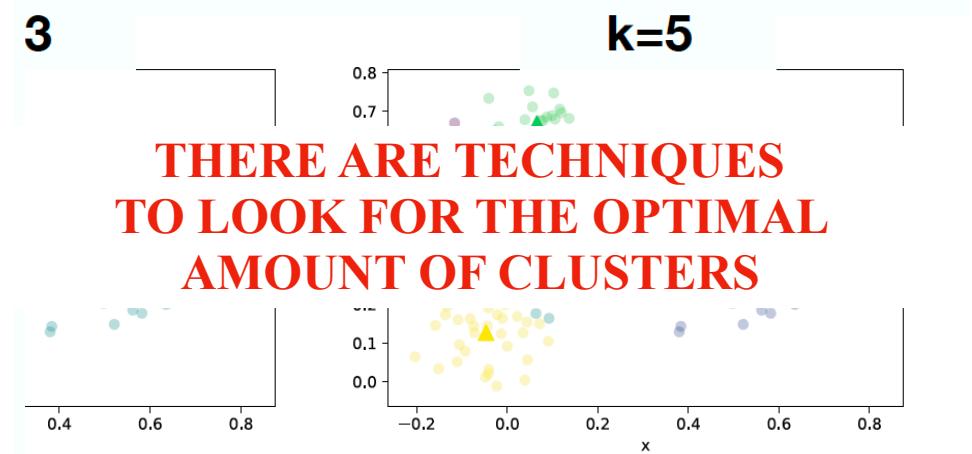
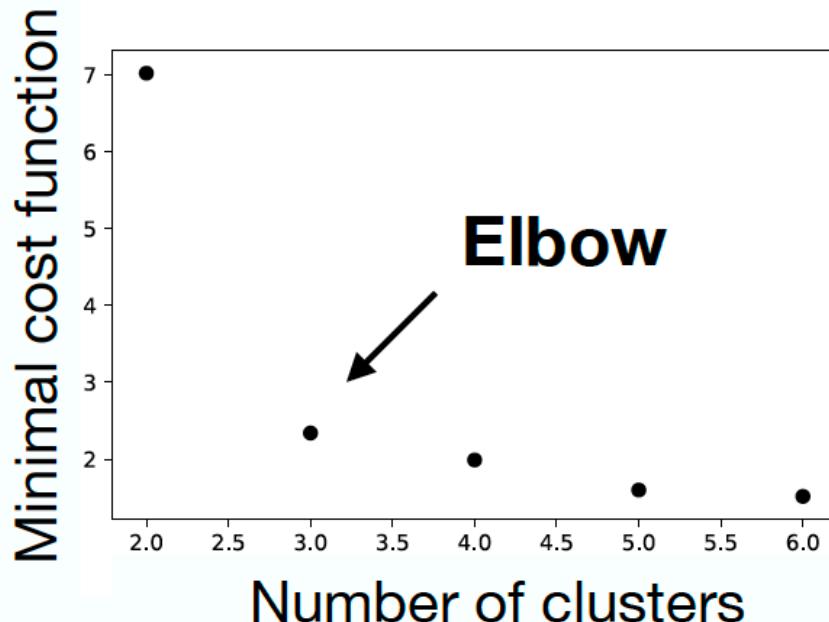


K-MEANS IS VERY SENSITIVE TO SEVERAL CRITICAL CHOICES

$$f(\vec{x}, \{a_1, a_2, \dots\}) = \vec{y}$$

Hyper-Parameters: The number of clusters needs to be

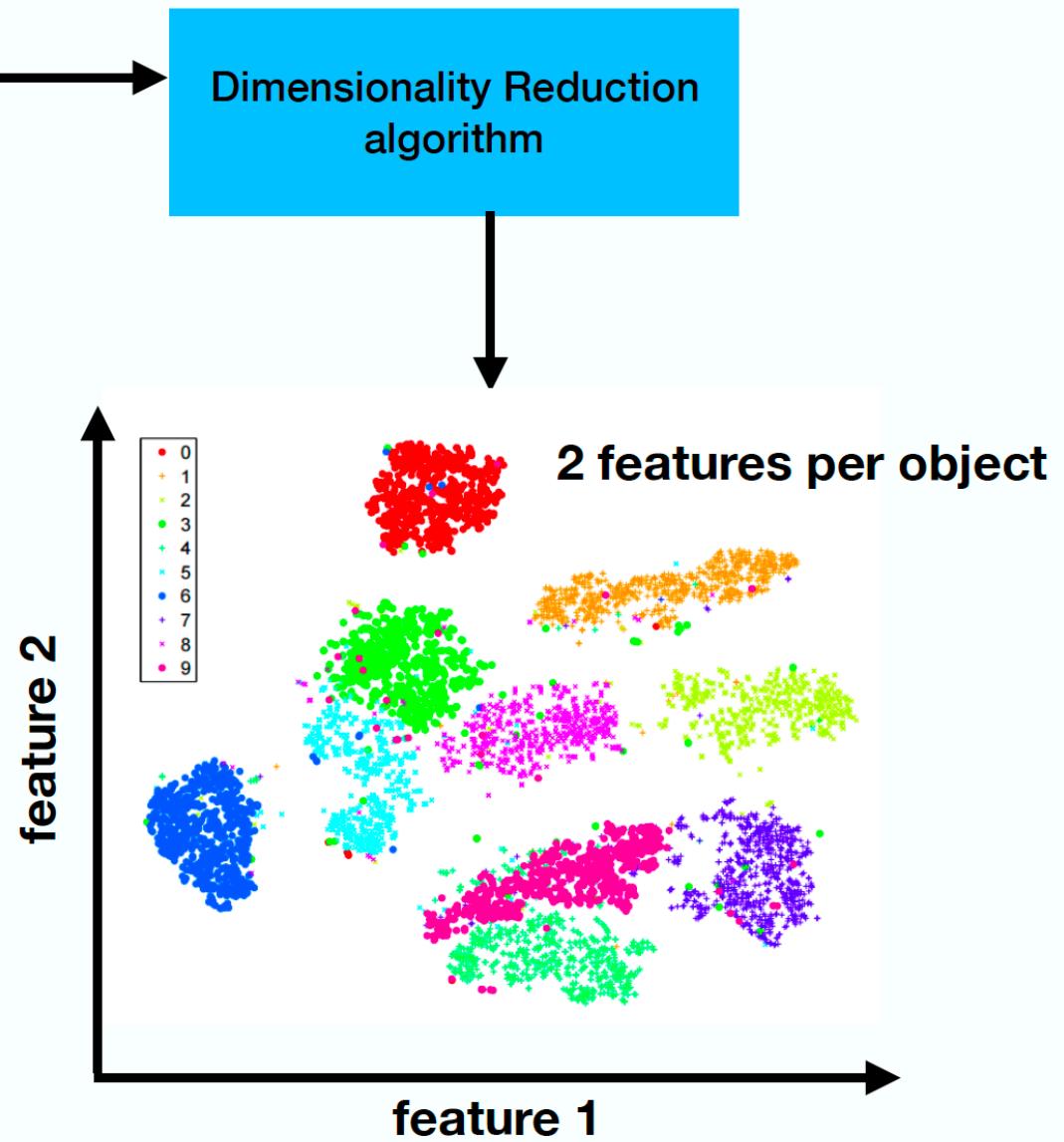
fived
98bis Boulevard Arago



3	6	0	3	0	1	1	3	9	3	1	5	0	4	9	6	8	7	1
0	5	6	9	8	8	4	1	4	4	4	6	4	5	3	3	4	3	4
0	4	3	7	7	5	0	5	4	2	0	9	8	1	2	4	9	3	5
1	1	1	7	4	7	7	2	6	5	1	8	2	4	1	1	5	6	5
7	0	9	5	6	3	2	6	6	7	1	5	2	3	2	3	5	6	
0	0	2	0	8	7	4	0	9	7	9	3	6	9	3	4	3	1	4
2	7	6	7	5	6	6	5	8	1	6	8	7	1	0	5	3	8	3
2	3	9	4	6	3	0	4	5	8	0	0	4	0	4	6	6	6	9
4	1	1	4	1	3	1	2	3	4	8	1	5	5	0	7	9	4	8

28 x 28 features per object

DIMENSIONALITY REDUCTION ALGORITHMS



WHY IS DIMENSIONALITY REDUCTION USEFUL?

- DATA COMPRESSION
- IMPROVE PERFORMANCE OF SUPERVISED LEARNING BY REDUCING REDUNDANCY IN ORIGINAL FEATURES
- DATA VISUALIZATION / INTERPRETATION
- UNCOVERING COMPLEX TRENDS
- OUTLIERS: UNKNOWN UNKNOWNS

THERE ARE TYPICALLY TWO TYPES OF DIMENSIONALITY REDUCTION ALGORITHMS

1. PROTOTYPE BASED

WE LOOK FOR PROTOTYPES THAN COMBINED, CAN REPRESENT MOST OF
THE VARIANCE CONTAINED IN THE DATA

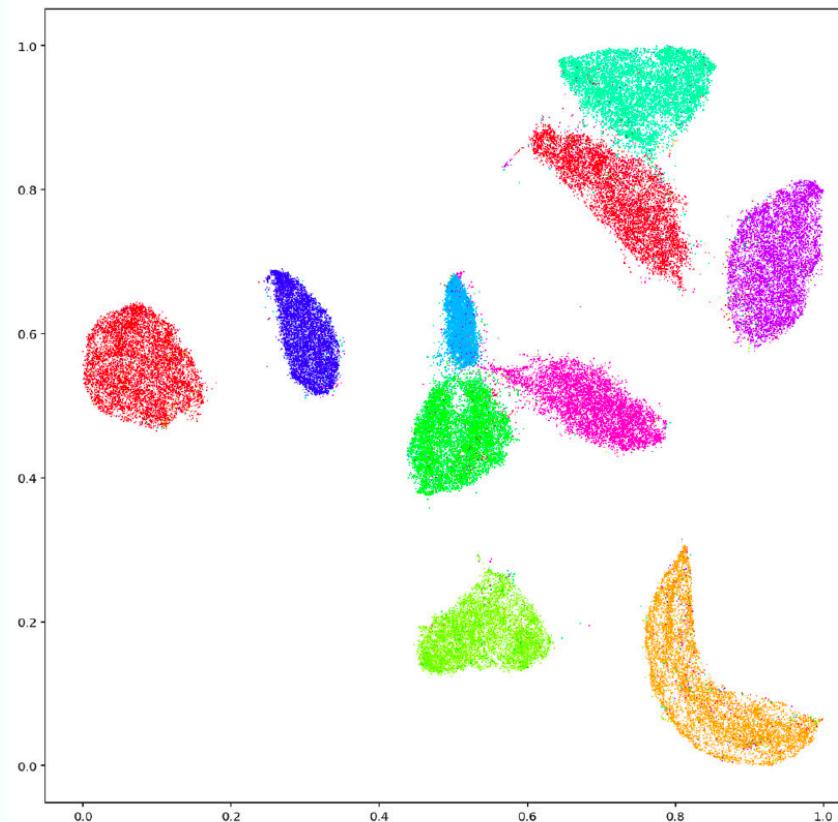


EXAMPLES ARE: PCA, SOMS, SVD, ICA...

THERE ARE TYPICALLY TWO TYPES OF DIMENSIONALITY REDUCTION ALGORITHMS

1. EMBEDDING

THEY EMBED A HIGH DIMENSIONAL DATASET INTO A LOWER DIMENSION

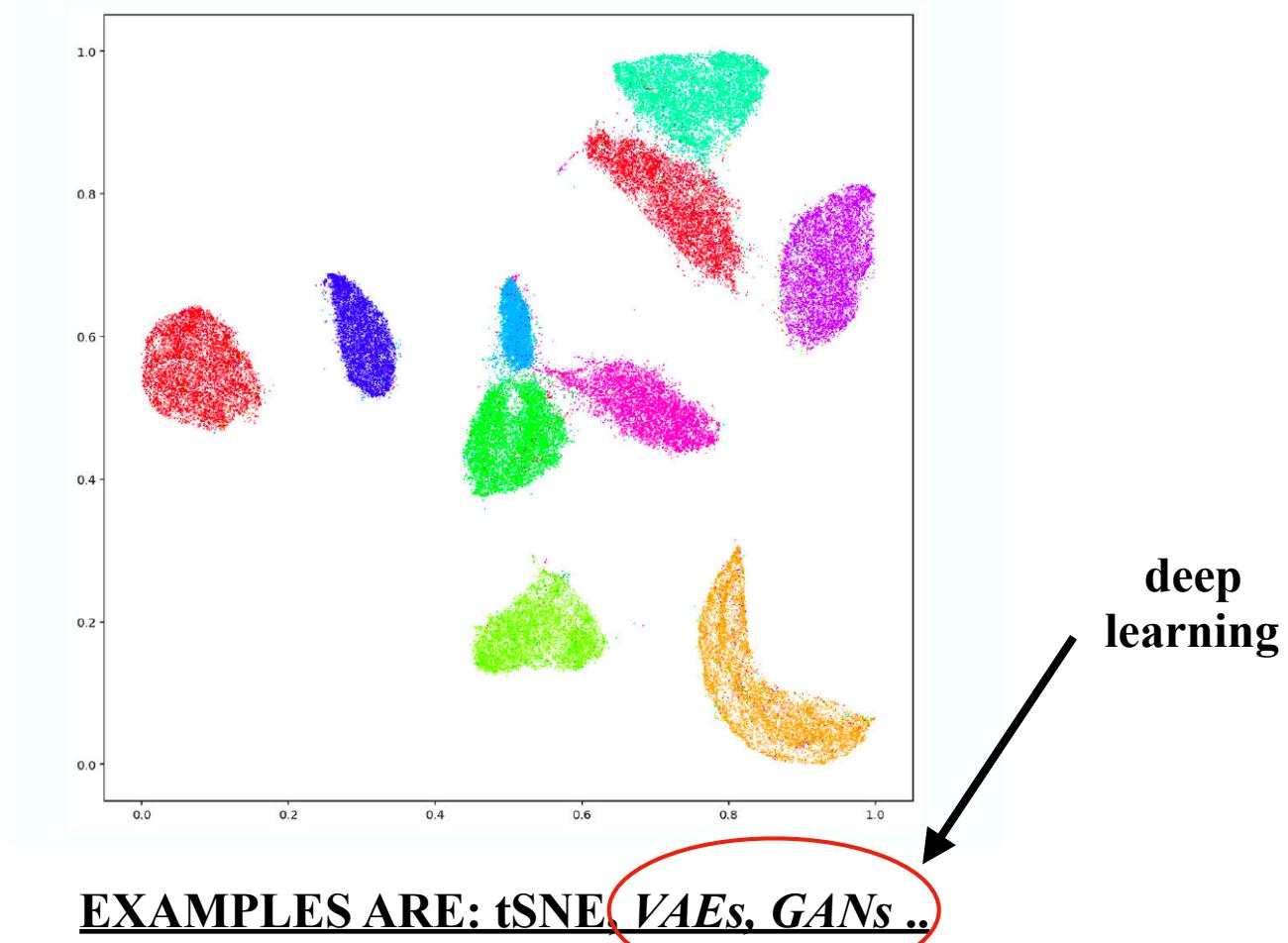


EXAMPLES ARE: tSNE, VAEs, GANs ..

THERE ARE TYPICALLY TWO TYPES OF DIMENSIONALITY REDUCTION ALGORITHMS

1. EMBEDDING

THEY EMBED A HIGH DIMENSIONAL DATASET INTO A LOWER DIMENSION

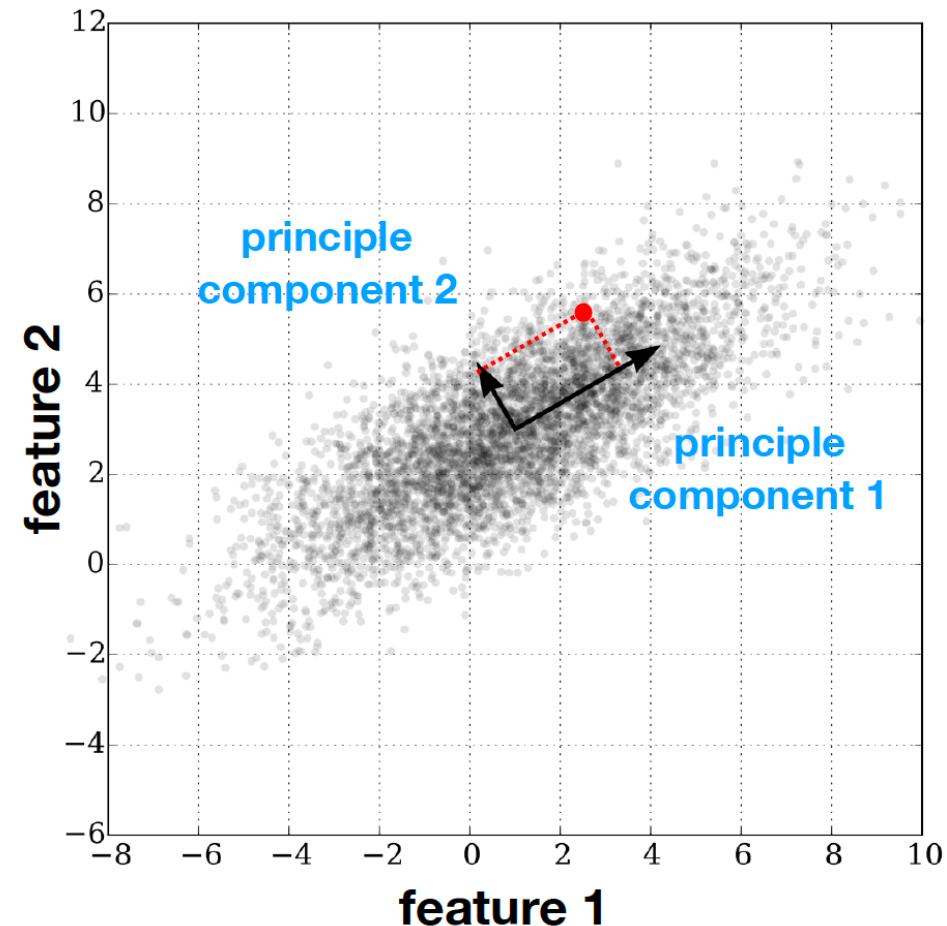


PRINCIPAL COMPONENT ANALYSIS

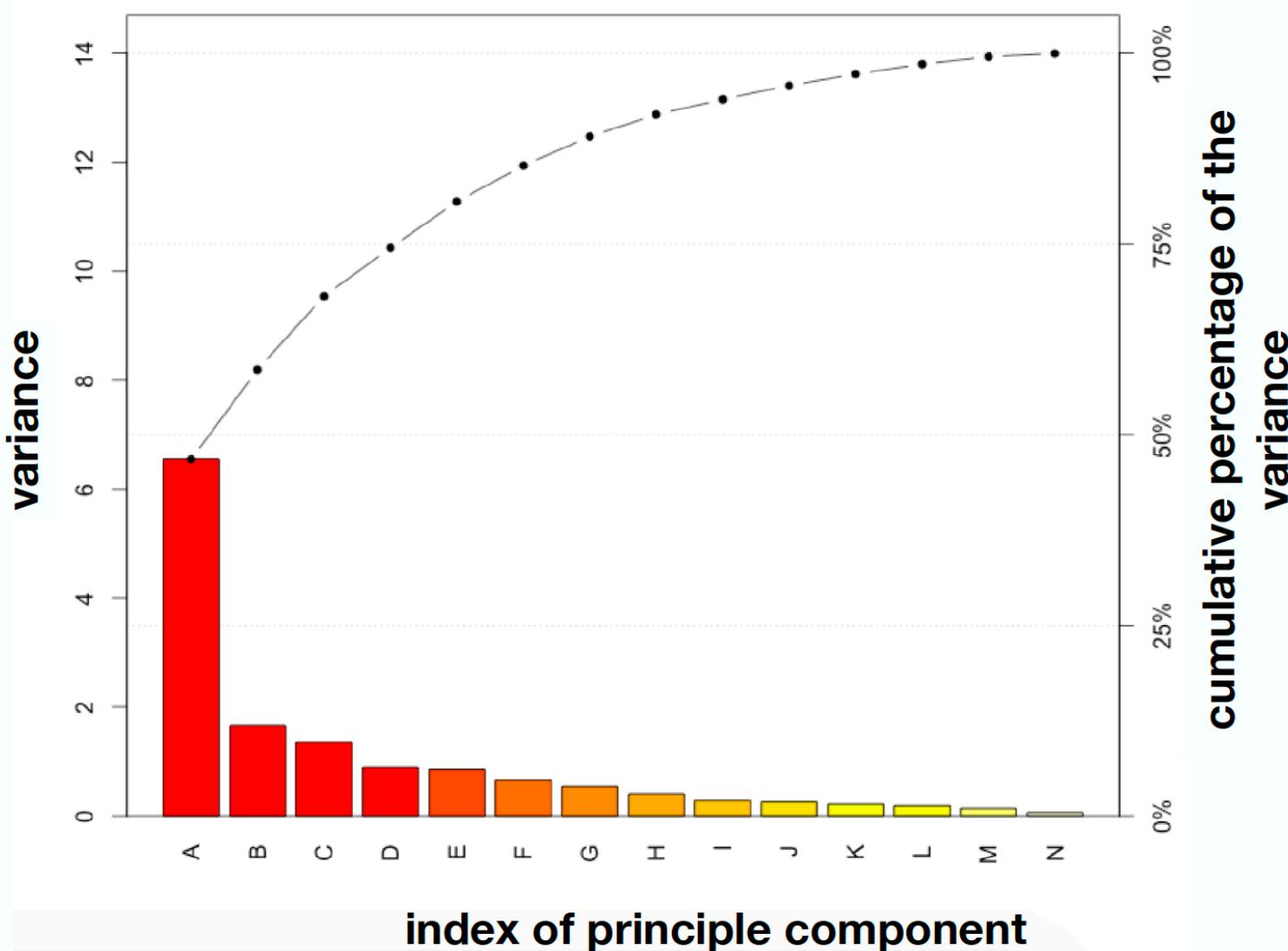
PCA CONVERT A SET OF
(CORRELATED) VARIABLES INTO A
SET
OF VALUES LINEARLY
UNCORRELATED

1. THE FIRST PRINCIPLE COMPONENT (“PROTOTYPE”), HAS THE LARGEST POSSIBLE VARIANCE

2. THE FOLLOWING COMPONENTS HAVE THE LARGEST VARIANCES WITH THE ADDITIONAL CONSTRAINT THAT THEY ARE ORTHOGONAL TO THE PRECEDING COMPONENTS



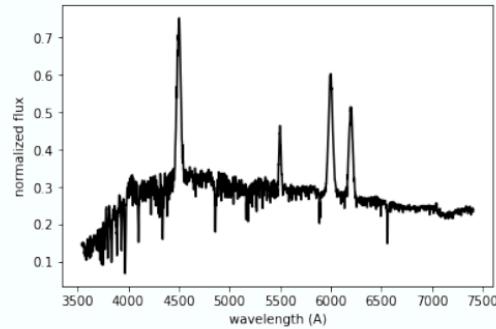
PRINCIPAL COMPONENT ANALYSIS



IT RESULTS IN DATA COMPRESSION,
BY REPRESENTING EACH OBJECT AS A PROJECTION OF THE FIRST PRINCIPLE COMPONENTS

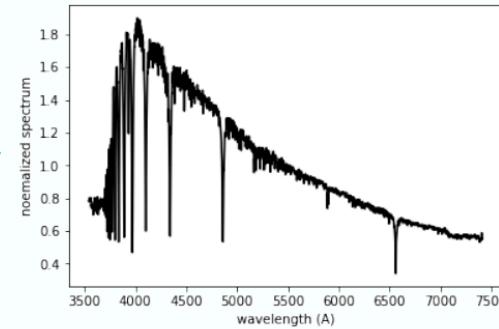
PRINCIPAL COMPONENT ANALYSIS

observed object



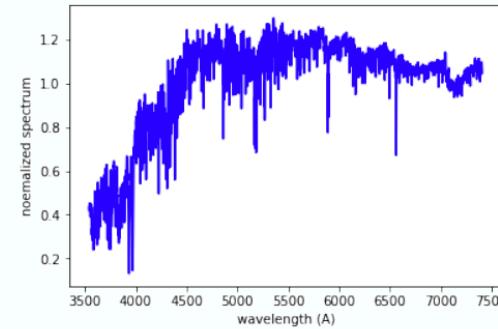
= A *

principle comp. 1



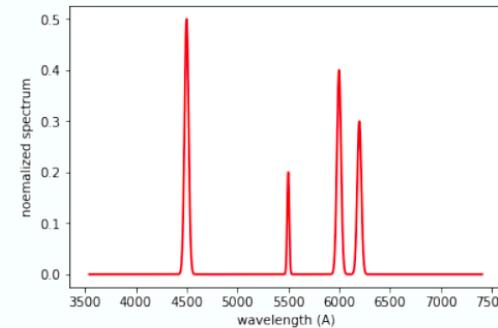
+ B *

principle comp. 2



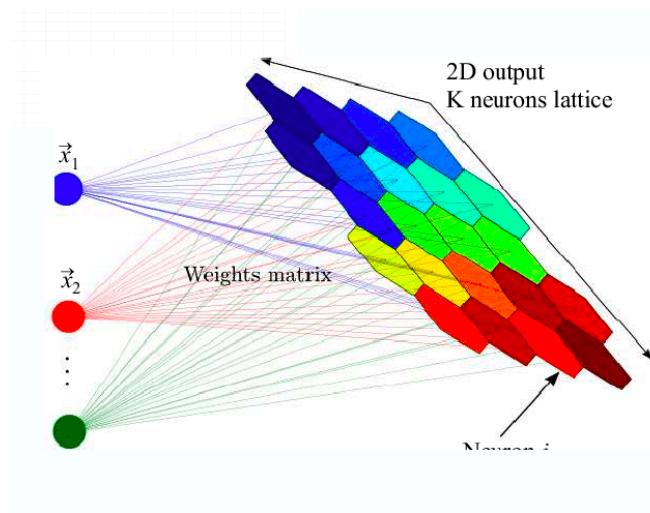
+ C *

principle comp. 3

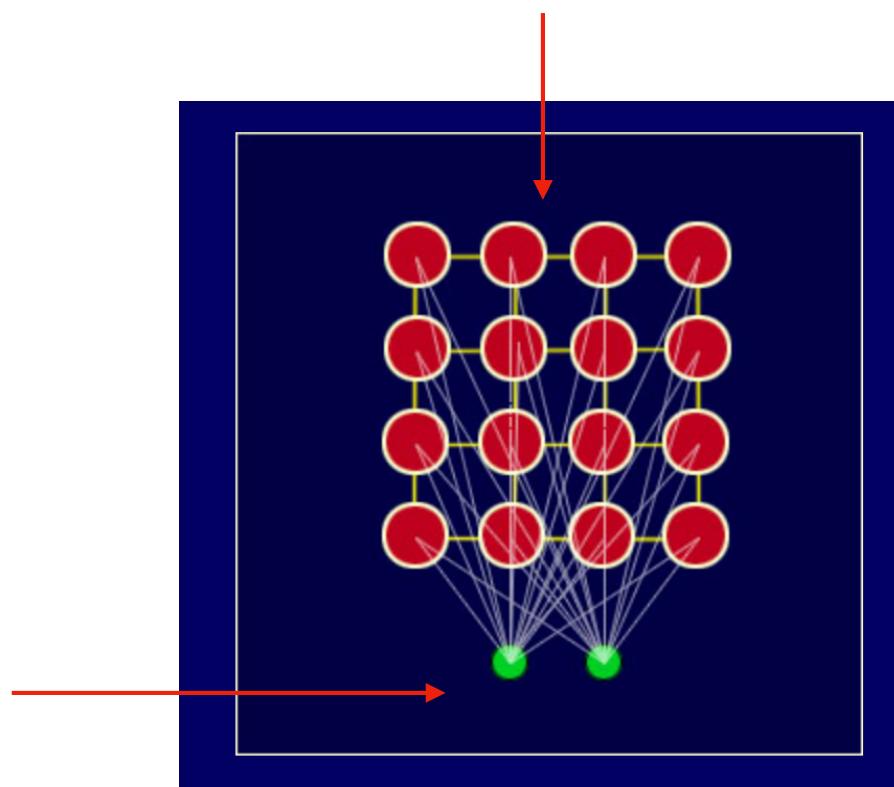


**OBJECTS ARE
DIVIDED INTO MAJOR
PROTOTYPES AND ALL
OBJECTS CAN BE
OBTAINED AS LINEAR
COMBINATIONS**

SELF ORGANIZING MAPS



**COMPLETION LAYER
("NEURONS" ARRANGE GEOMETRICALLY)**



EACH "NODE" IN THE SOM HAS THE SAME NUMBER OF WEIGHTS THAN INPUT VECTORS

SELF ORGANIZING MAPS

- 1. EACH NODE'S WEIGHTS ARE INITIALIZED RANDOMLY**

SELF ORGANIZING MAPS

2. A VECTOR IS CHOSEN AT RANDOM FROM THE SET OF TRAINING DATA AND PRESENTED TO THE LATTICE

$$D = \sqrt{\sum_{i=0}^{i=n} (V_i - W_i)^2}$$

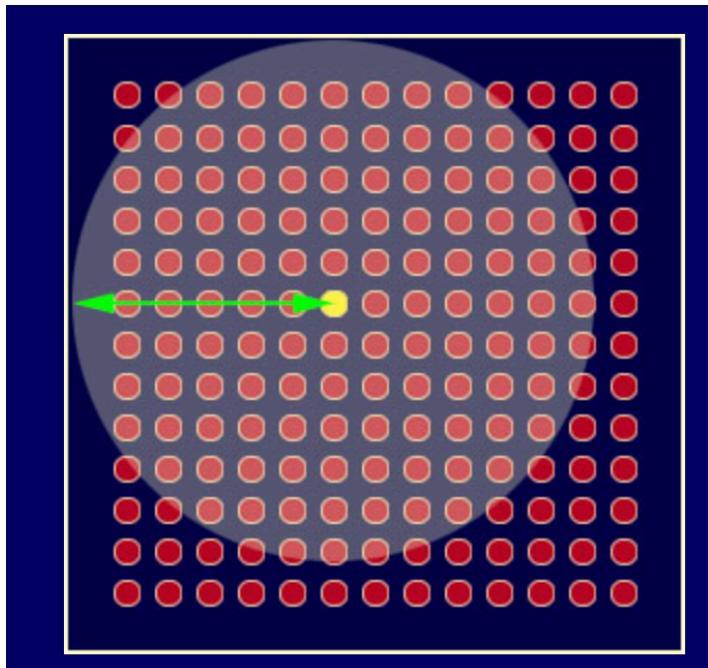
3. EVERY NODE IS EXAMINED TO CALCULATE WHICH ONE'S WEIGHTS ARE MOST LIKE THE INPUT VECTOR. THE WINNING NODE IS COMMONLY KNOWN AS THE BEST MATCHING UNIT (BMU).

SELF ORGANIZING MAPS

4. THE RADIUS OF THE NEIGHBORHOOD OF THE BMU IS NOW CALCULATED. THIS IS A VALUE THAT STARTS LARGE, TYPICALLY SET TO THE 'RADIUS' OF THE LATTICE, BUT DIMINISHES EACH TIME-STEP.

$$\sigma(t) = \sigma_0 \times e^{-t/\lambda}$$

THE RADIUS
DECREASES OVER TIME



SELF ORGANIZING MAPS

5. EACH NEIGHBORING NODE'S (THE NODES FOUND IN STEP 4) WEIGHTS ARE ADJUSTED TO MAKE THEM MORE LIKE THE INPUT VECTOR. THE CLOSER A NODE IS TO THE BMU, THE MORE ITS WEIGHTS GET ALTERED.

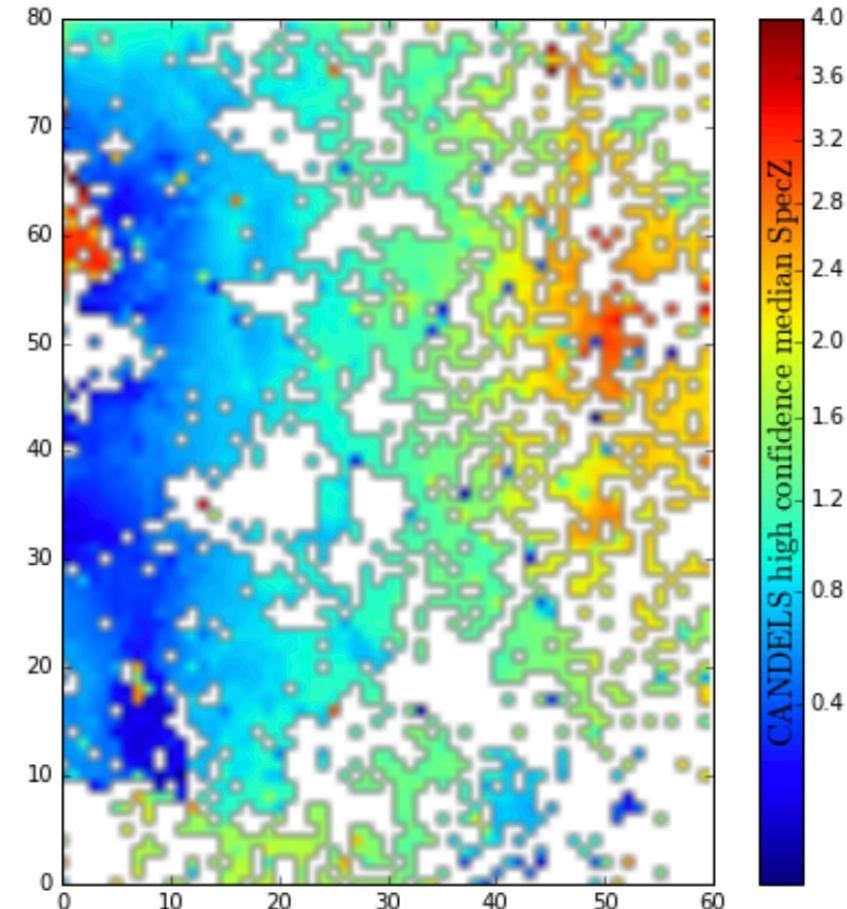
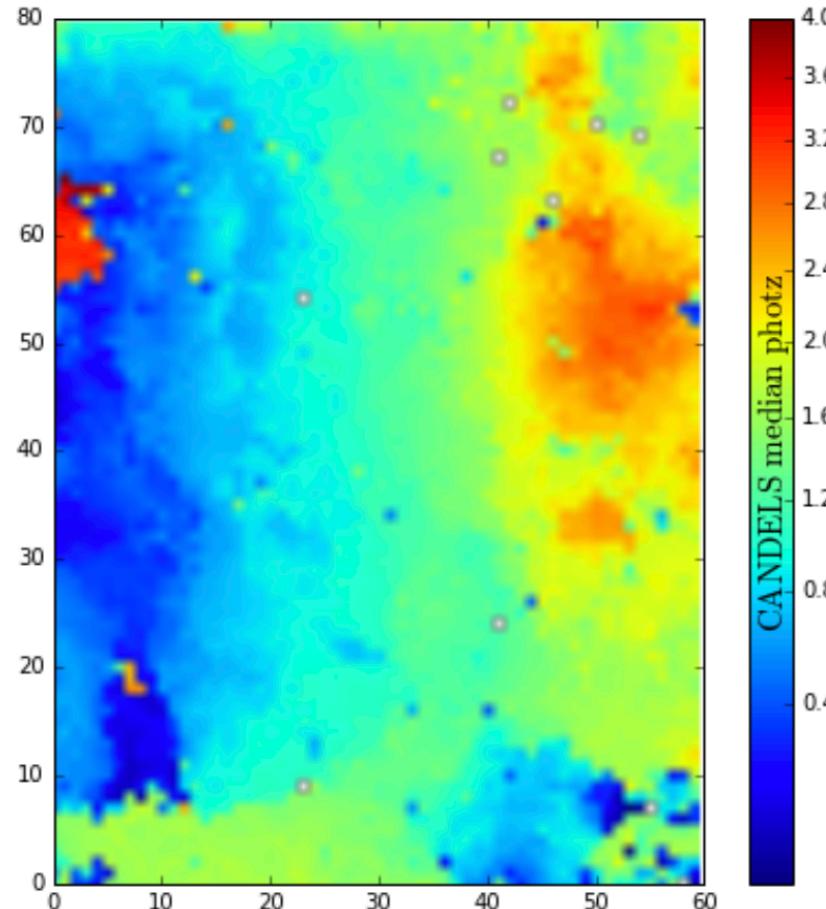
$$W(t + 1) = W(t) + \Psi(t)\lambda(t)[V(t) - W(t)]$$

The diagram illustrates the components of the weight update equation for a Self-Organizing Map. It shows the equation $W(t + 1) = W(t) + \Psi(t)\lambda(t)[V(t) - W(t)]$ with arrows pointing from each term to its corresponding component:

- A diagonal arrow points from the term $\Psi(t)\lambda(t)$ to the label **CORRECTION FACTOR**.
- A vertical arrow points from the term $V(t)$ to the label **FEATURES**.
- A vertical arrow points from the term $W(t)$ to the label **WEIGHTS**.
- A horizontal arrow points from the term $[V(t) - W(t)]$ to the label **LEARNING RATE**.

SELF ORGANIZING MAPS

6. REPEAT FROM 2 FOR N ITERATIONS



HEMMATI ET AL. 2018

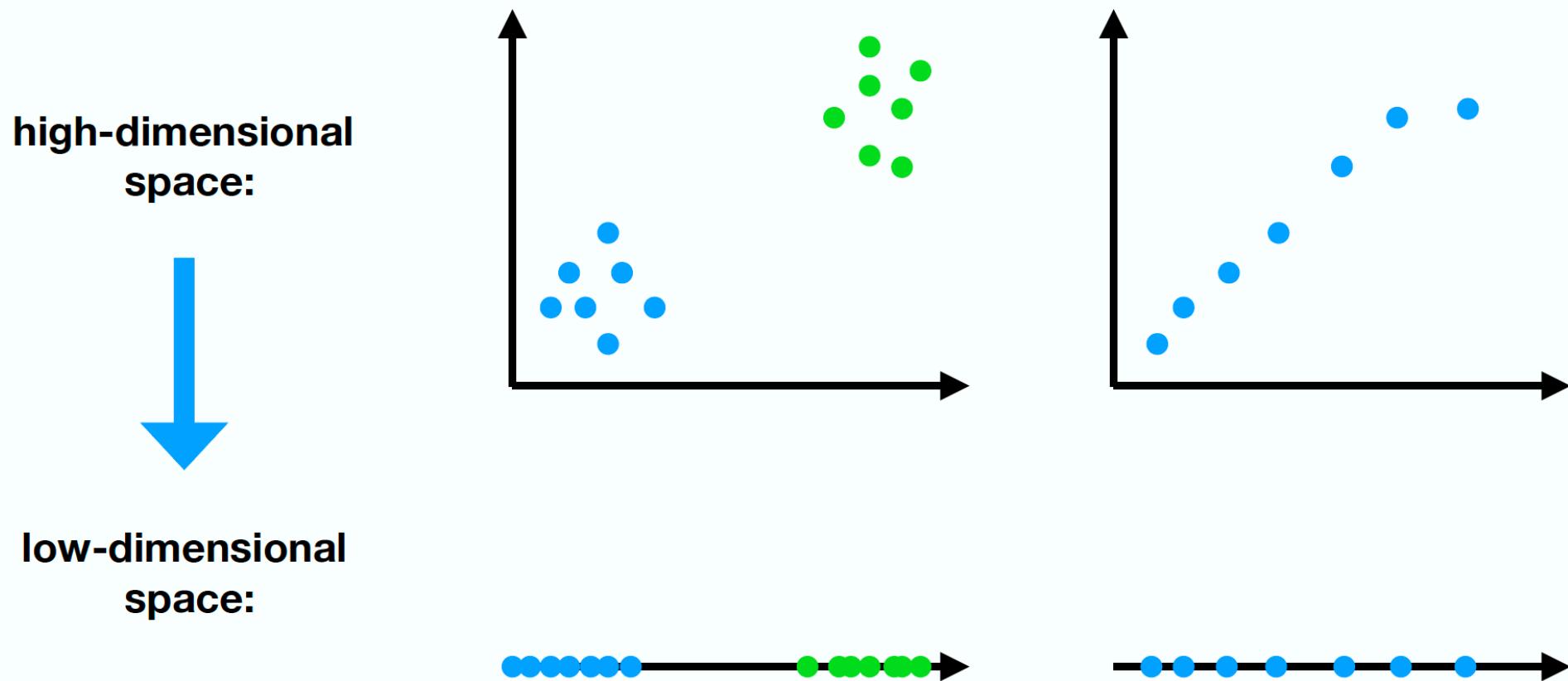
SOMs FOR PHOTOMETRIC REDSHIFTS

tSNE

t-distributed stochastic neighbor embedding

Embedding high-dimensional data in a low dimensional space (2 or 3)

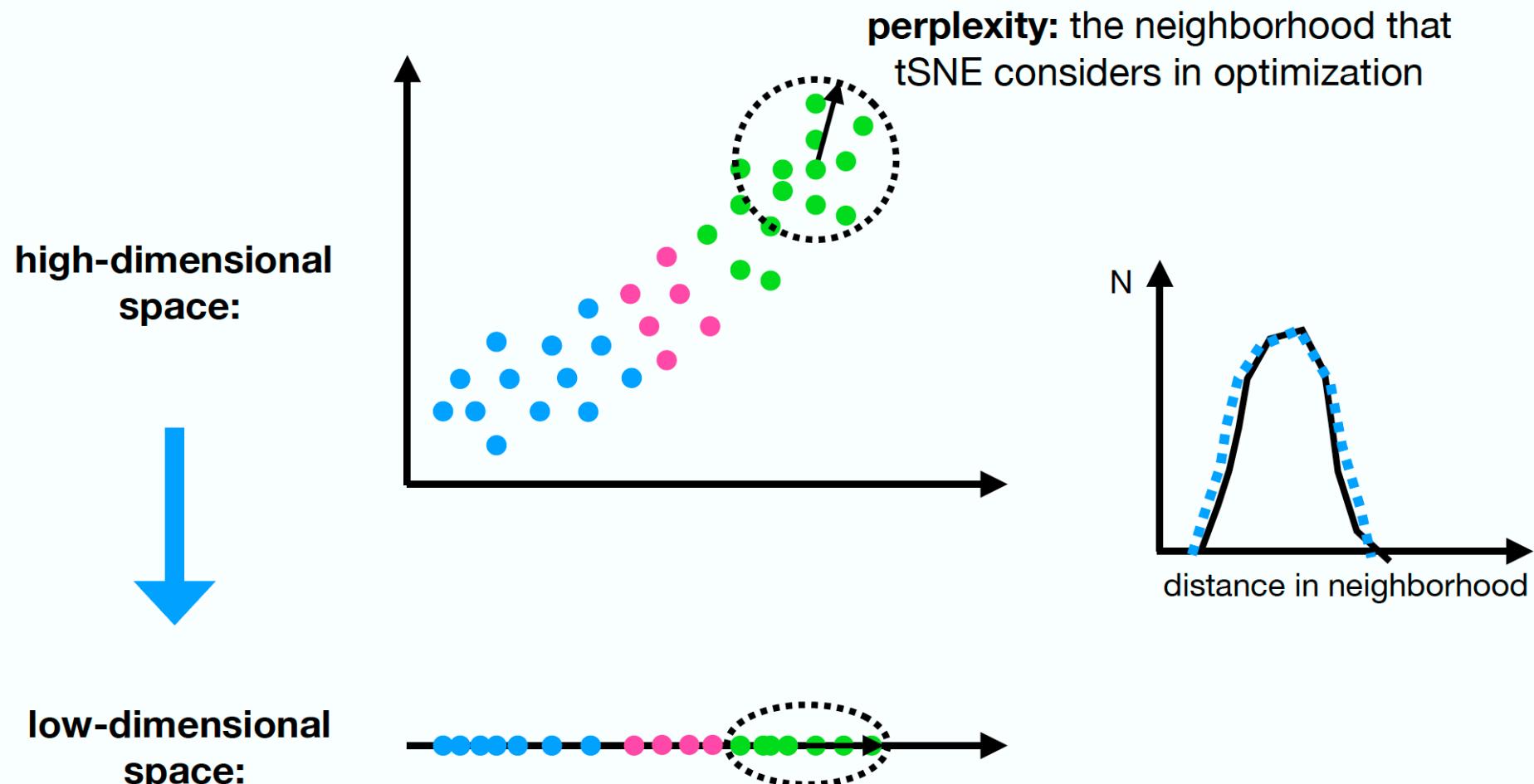
Input: (1) raw data, extracted features, or a distance matrix
(2) hyper-parameters: **perplexity**



tSNE

t-distributed stochastic neighbor embedding

Intuition: tSNE tries to find a low-dimensional embedding that preserves, as much as possible, the [distribution of distances](#) between different objects.



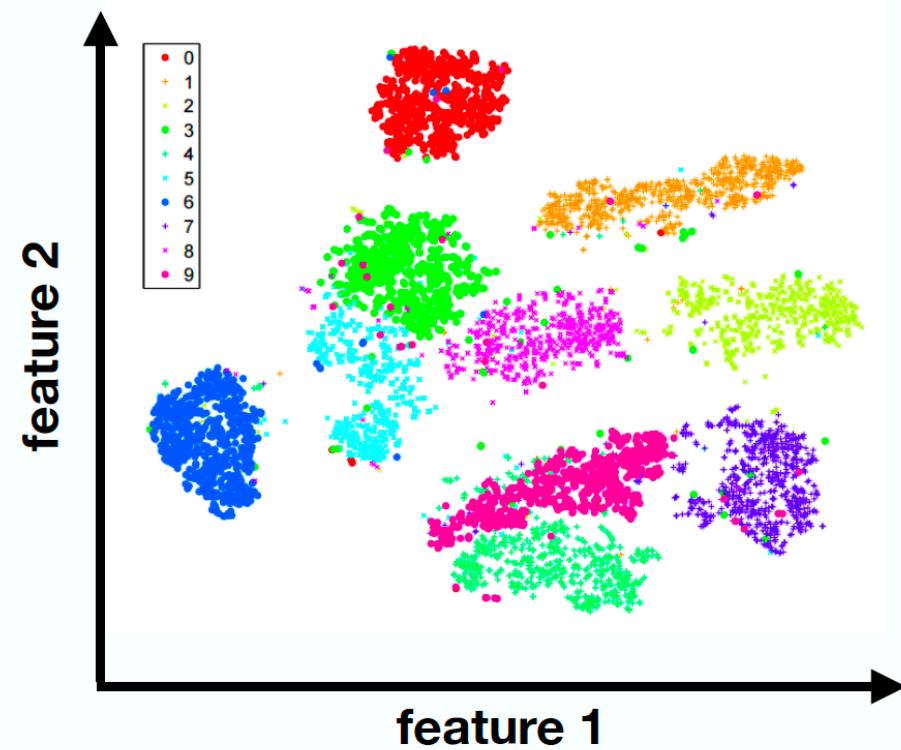
Credit: Baron IAC XXX Winter

tSNE

t-distributed stochastic neighbor embedding

3 6 0 3 0 / 1 3 9 3 1 5 0 4 9 6 8 7 1
0 5 6 9 8 8 4 1 4 4 4 6 4 5 3 3 4 3 4
0 4 3 7 7 5 0 5 4 2 0 9 8 1 2 4 9 3 5
1 1 1 7 4 7 7 2 6 5 1 8 2 4 1 1 5 6 5
7 0 9 5 6 3 2 6 6 7 1 5 2 3 2 3 5 6
0 0 2 0 8 7 4 0 9 7 9 3 6 9 3 4 3 1 4
2 7 6 7 5 6 6 5 8 1 6 8 7 1 0 5 3 8 3
2 3 9 4 3 0 4 5 8 0 0 4 0 4 6 6 6 9 3
4 1 1 4 1 3 1 2 3 4 8 1 5 5 0 7 9 4 8

28 x 28 features per object



Credit: Baron IAC XXX Winter