

S&DS 365 / 565  
Data Mining and Machine Learning

shift from  
supervised learning to

# Unsupervised Learning

Yale

Loose idea

Given unlabeled data, find structure.

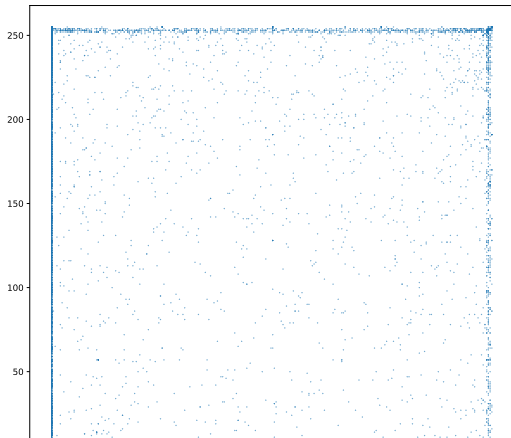
Methods { Clustering  
t-SNE  
Bayesian

Topic Modelling of docs, e-mails

Given unlabeled data, find structure.

2-D  
( $x_i, y_i$ )

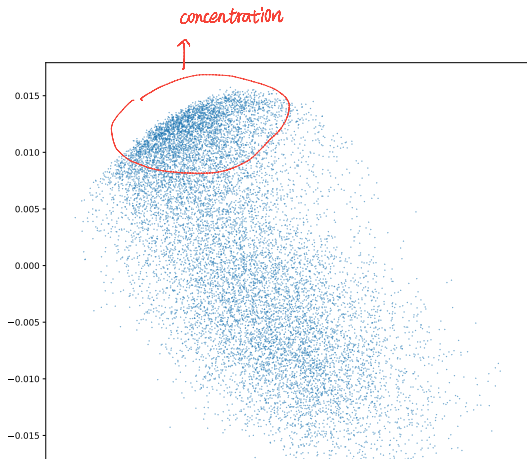
Suppose I have a dataset that is  $12593 \times 784$ . How can I visualize it?  
Maybe plot some coordinates against each other?



Given unlabeled data, find structure.

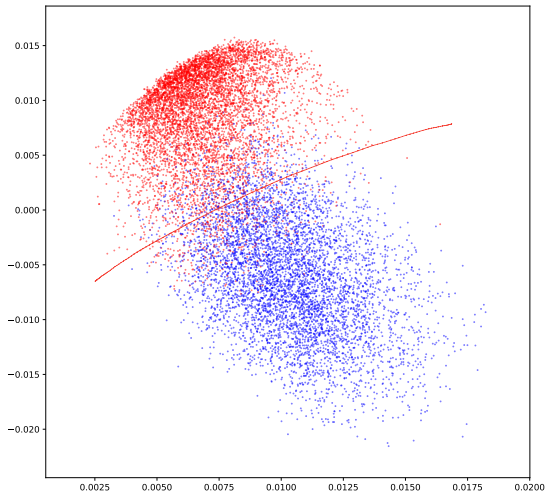
*baseline*

PCA is a workhorse method for computing interesting directions of the data. Here I use PCA to plot “interesting” direction of the data. We will discuss what “interesting” means later.



# Let's add some color.

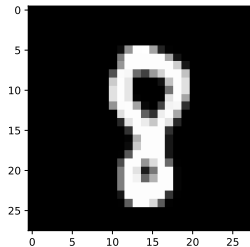
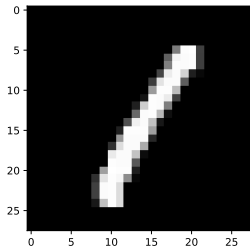
*based on the label*



Actually, I generated the plots from data with labels. Images to be exact, the pattern will become very clear to you.

Actually, I generated the plots from data with labels. Images to be exact, the pattern will become very clear to you.

*2 labels*



Now: clustering

Aim: Find good representatives

Example: Identify distinct communities of butterflies based on wing size, mass, color

*features* (under "size, mass, color")  
*label* (under "communities")

Example: Digits images, find representatives (obviously, but since we have labels we can check our clustering)



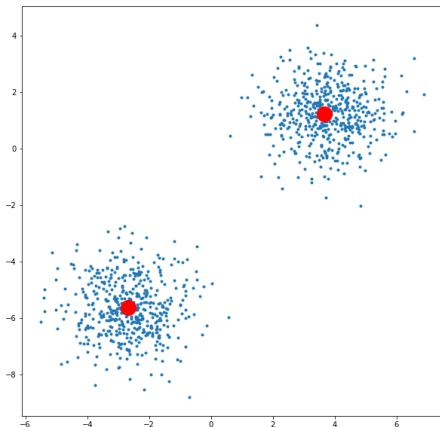
Clustering is a loose concept

Start with **k-means** clustering (explored in previous HW)<sup>1</sup>

Dataset:  $x_i \in \mathbb{R}^p$   $i \in [n]$

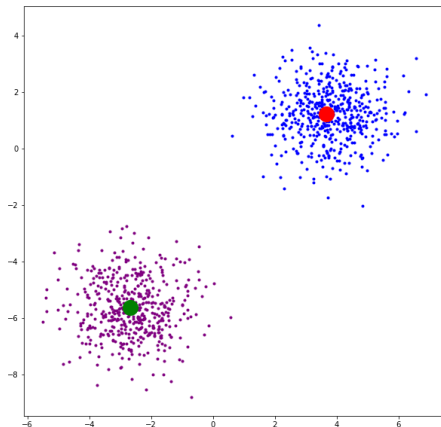
Goal: find **vectors (centers)**  $\mu_j$   $j \in [K]$  such that they represent the data well.

$k=2$  plot in 2-D



We can assign points to each center.

↑  
mapping  $\pi(i) \in [2]$



We assign purple points to center green and blue points to center red

This provides a mapping  $\pi$  of points to centers. So  $\pi(i) \in [K]$ . There are  $K$  different centers, so that explains the  $K$  part of  $K$ -means.

To understand the means part we can now present how the  $K$ -means algorithm decides if centers represent the data “well.”

This provides a mapping  $\pi$  of points to centers. So  $\pi(i) \in [K]$ . There are  $K$  different centers, so that explains the  $K$  part of  $K$ -means.

To understand the means part we can now present how the  $K$ -means algorithm decides if centers represent the data “well.”

Goal: Find centers  $\overset{\text{representative}}{\mu_j}$  such that

$$\mu_j = \arg \min_{\mu_{\pi(i)}} \sum_{i=1}^n \|x_i - \mu_{\pi(i)}\|_2^2$$

is minimized. Here,  $\pi(i)$  is actually a function that maps example  $i$  to centers  $j$ . So we have to find both  $\mu$  and  $\pi$ .

In general the computation is intractable. However, if we are given  $\pi$  finding  $\mu_j$  is easy. Homework assignment to check to see that

$$\text{take mean} \quad \hat{\mu}_j = \frac{1}{n_j} \sum_{i|\pi(i)=j} x_i \quad j \in [k]$$

where  $n_j = \sum_{i|\pi(i)=j} 1$ . This explains the **means** part of  $K$ -means!  
*number of constituents in group j*

*map example  $i$  to the closest representative  $\mu_j$*

Similarly, given  $\mu_j$  finding  $\pi$  is also very easy.

$$\pi(i) = \arg \min_j \|x_i - \mu_j\|_2^2$$

These observations yield a natural **alternating minimization** algorithm.

*k-means clustering algorithm*  
*random*

Starting with initial guess for  $\mu_j$

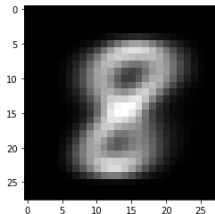
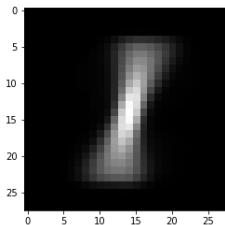
- Compute  $\pi$
- Then compute  $\mu_j$
- repeat until convergence



Let's try it out on the two digits with  $K = 2$ .

fuzzy images coz people write digits "1" "8" differently  
The images capture all possible 笔迹 then average them

output 2 clusters :



Let's try it out on the two digits with  $K = 3$ .

