

# Ch 12.1, 12.4: Unsupervised Learning & Clustering

Lecture 32 - CMSE 381

Prof. Elizabeth Munch

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Mon, Nov 25, 2024

# Announcements

## Last time:

- Convolutional Neural Nets

## This lecture:

- Clustering (Just hierarchical clustering)

## Announcements:

- No more homework!
- Weds: Project office hours, zoom only, send a message on slack!
- Mon Dec 2: Review - Bring questions!
- Weds Dec 4: Exam
  - ▶ Content since 2nd Exam (Ch 7 and on)
  - ▶ One page (8.5x11) handwritten cheat sheet
  - ▶ Calculator if you want it

Lec #	Date		Reading	HW
21	Mon 10/28	Polynomial & Step Functions	7.1,7.2	
22	Wed 10/30	Step Functions; Basis functions; Start Splines	7.2 - 7.4	
23	Fri 11/1	Regression Splines	7.4	HW #6 Due
24	Mon 11/4	Decision Trees	8.1	Sun 11/3
25	Wed 11/6	Class Cancelled (Dr Munch out of town)		
26	Fri 11/8	Random Forests	8.2.1, 8.2.2	HW #7 Due
27	Mon 11/11	Maximal Margin Classifier	9.1	Sun 11/10
28	Wed 11/13	SVC	9.2	
29	Fri 11/15	SVM	9.3, 9.4	HW #8 Due
30	Mon 11/18	Single layer NN	10.1	Sun 11/17
31	Wed 11/20	Multi Layer NN	10.2	
32	Fri 11/22	CNN	10.3	HW #11
33	Mon 11/25	TBD: Unsupervised learning/clustering	12.1, 12.4?	Due Sun 11/24
	Wed 11/27	Virtual: Project office hours		
	Fri 11/29	No class - Thanksgiving		
	Mon 12/2	Review		
	Wed 12/4	Midterm #3		
	Fri 12/6	No class - EGR Design Day		Project due

# Section 1

## Unsupervised learning

# Supervised vs Unsupervised Learning

**Supervised**

**Unsupervised**

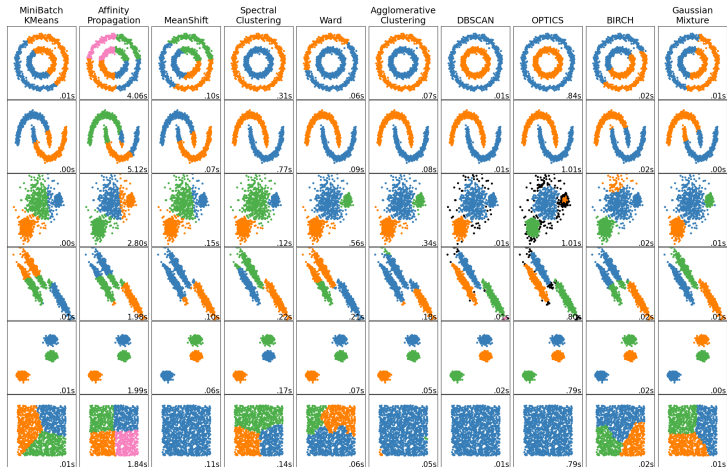
## Some examples of unsupervised problems

- Assay gene expression levels in 100 patients with breast cancer, looking for subgroups with similar qualities
- Online shopping: find groups of shoppers with similar browsing and purchase histories and show relevant related products.
- Search engine picking results to show

## Section 2

### Clustering

# Big idea

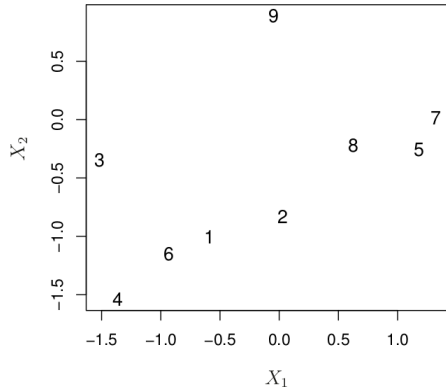
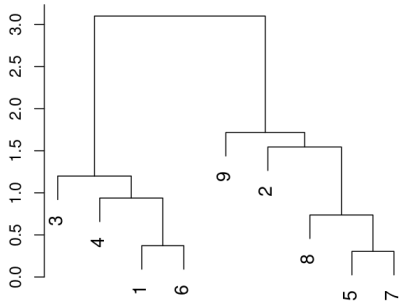


## Section 3

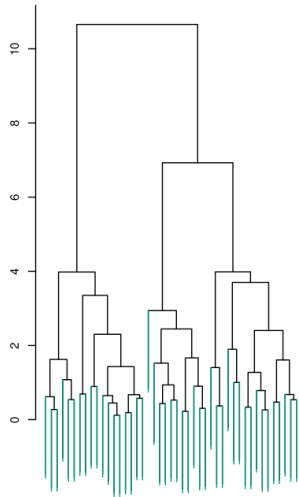
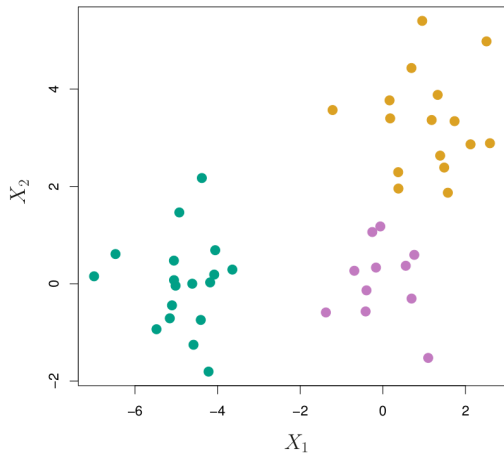
# Hierarchical Clustering



# Dendrogram



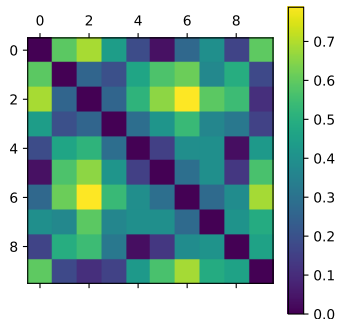
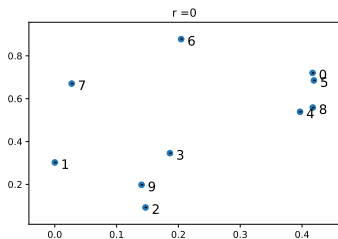
# A bigger example



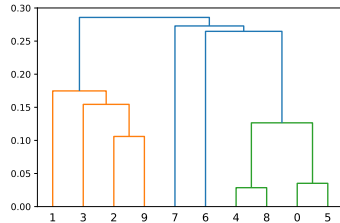
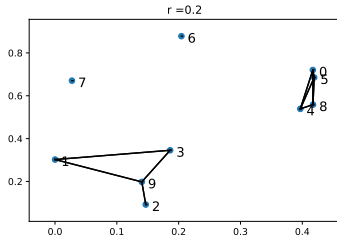
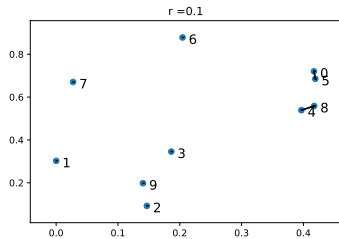
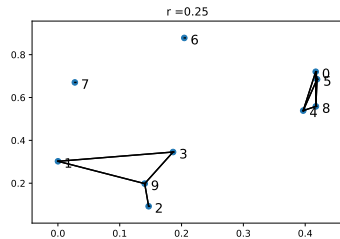
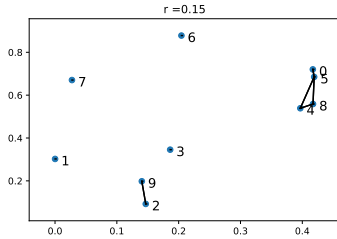
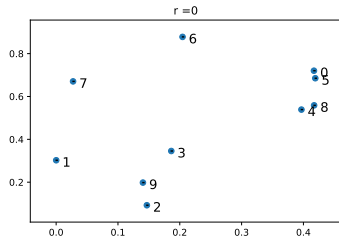
# Single linkage

Distance between cluster  $A$  and cluster  $B$ :  
Smallest distance between the points

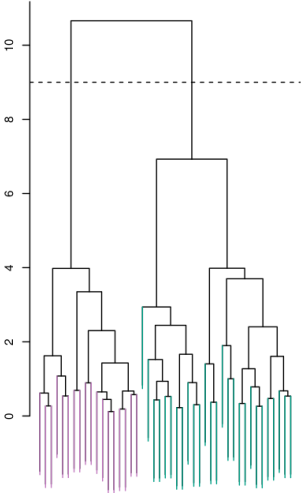
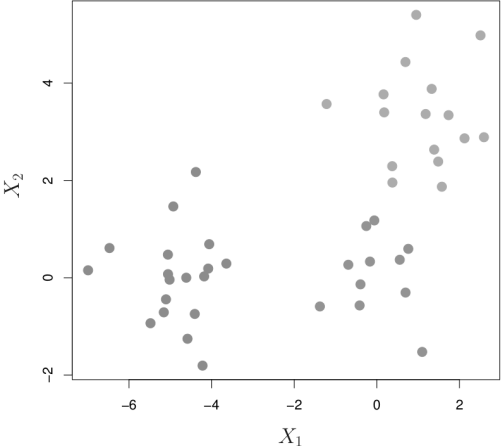
$$L(A, B) = \min_{a \in A, b \in B} \|a - b\|$$



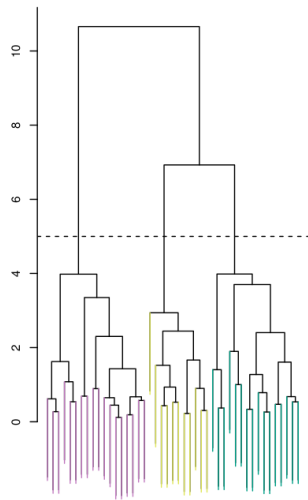
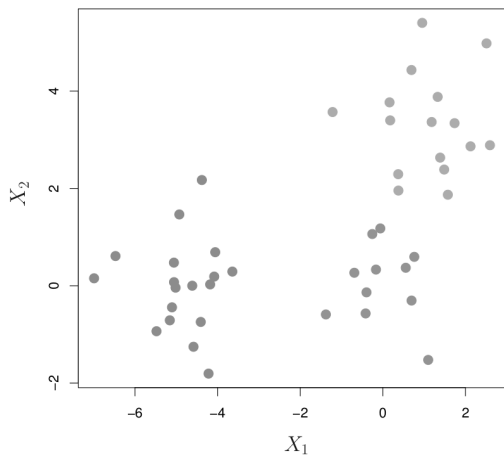
# Building the dendrogram



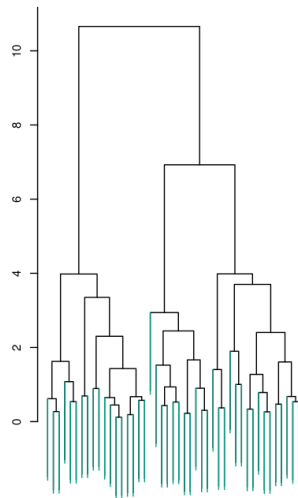
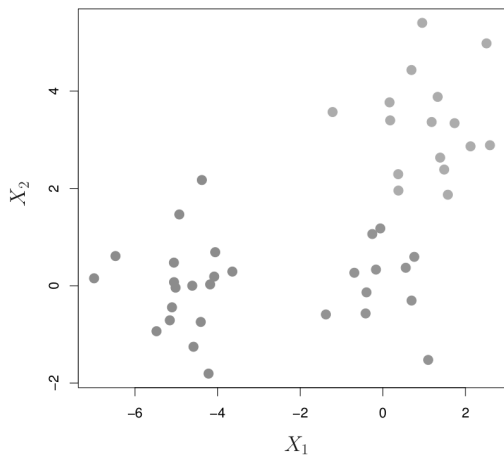
# How to get clusters



# How to get different clusters



# Can get any number of clusters

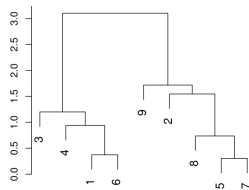


# Linkage

<i>Linkage</i>	<i>Description</i>
Complete	Maximal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>largest</i> of these dissimilarities.
Single	Minimal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>smallest</i> of these dissimilarities. Single linkage can result in extended, trailing clusters in which single observations are fused one-at-a-time.
Average	Mean intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>average</i> of these dissimilarities.
Centroid	Dissimilarity between the centroid for cluster A (a mean vector of length $p$ ) and the centroid for cluster B. Centroid linkage can result in undesirable <i>inversions</i> .

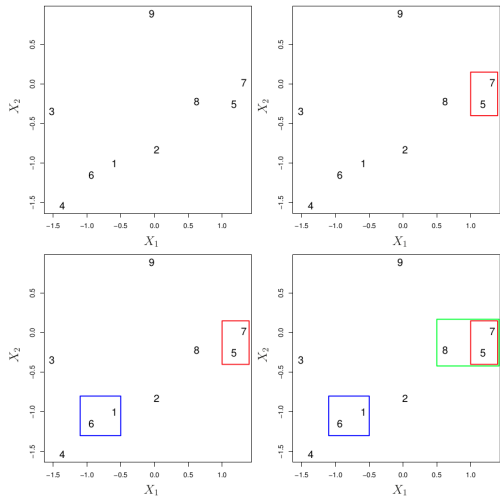


# Example with complete linkage



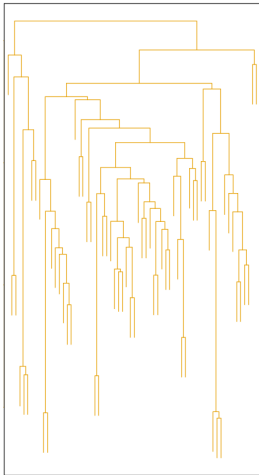
Distance between cluster  $A$  and cluster  $B$ :  
**Largest** distance between the points

$$L(A, B) = \max_{a \in A, b \in B} \|a - b\|$$

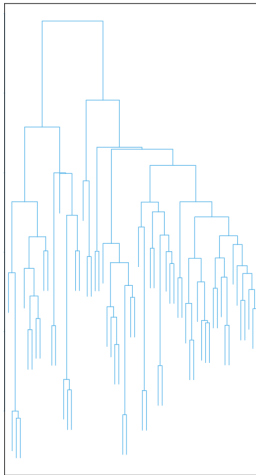


# Examples of different linkage

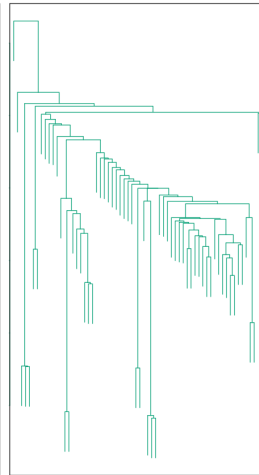
Average Linkage



Complete Linkage



Single Linkage



# Dependence on dissimilarity measure

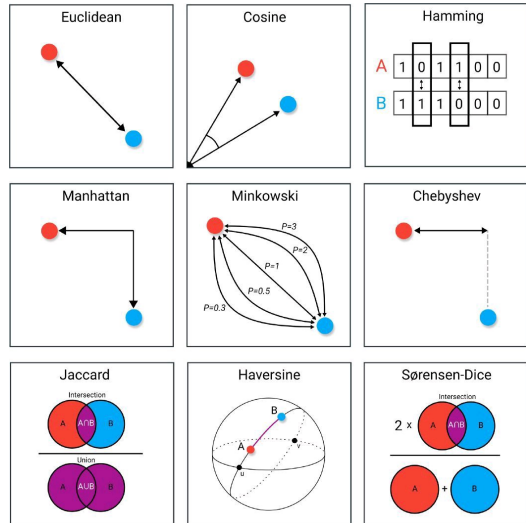


Photo Credit Link



# Next time

Lec #	Date		Reading	HW
21	Mon 10/28	Polynomial & Step Functions	7.1,7.2	
22	Wed 10/30	Step Functions; Basis functions; Start Splines	7.2 - 7.4	
23	Fri 11/1	Regression Splines	7.4	HW #6 Due
24	Mon 11/4	Decision Trees	8.1	Sun 11/3
25	Wed 11/6	Class Cancelled (Dr Munch out of town)		
26	Fri 11/8	Random Forests	8.2.1, 8.2.2	HW #7 Due
27	Mon 11/11	Maximal Margin Classifier	9.1	Sun 11/10
28	Wed 11/13	SVC	9.2	
29	Fri 11/15	SVM	9.3, 9.4	HW #8 Due
30	Mon 11/18	Single layer NN	10.1	Sun 11/17
31	Wed 11/20	Multi Layer NN	10.2	
32	Fri 11/22	CNN	10.3	HW #11
33	Mon 11/25	TBD: Unsupervised learning/clustering	12.1, 12.4?	Due Sun 11/24
	Wed 11/27	Virtual: Project office hours		
	Fri 11/29	No class - Thanksgiving		
	Mon 12/2	Review		
	Wed 12/4	<b>Midterm #3</b>		
	Fri 12/6	No class - EGR Design Day		Project due