

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023



HVERFORD
COLLEGE

Admin

- Lab 6 graded
- Presentation schedule up
- Lab today
 - Final project check-ins with all groups
 - Try to come to the same lab as your partner
- Candidate talk at 4:15pm TODAY
 - Tea at 4pm
 - Student lunch Wednesday 12:30-1:30pm

Outline for November 28

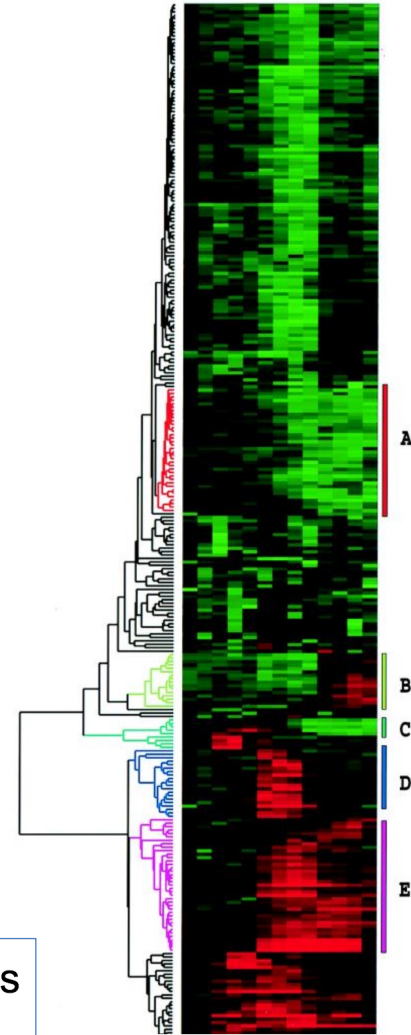
- Clustering overview
- K-means
- Gaussian Mixture Models (GMMs)

Outline for November 28

- Clustering overview
- K-means
- Gaussian Mixture Models (GMMs)

Applications of clustering

- Cluster genes with similar expression patterns

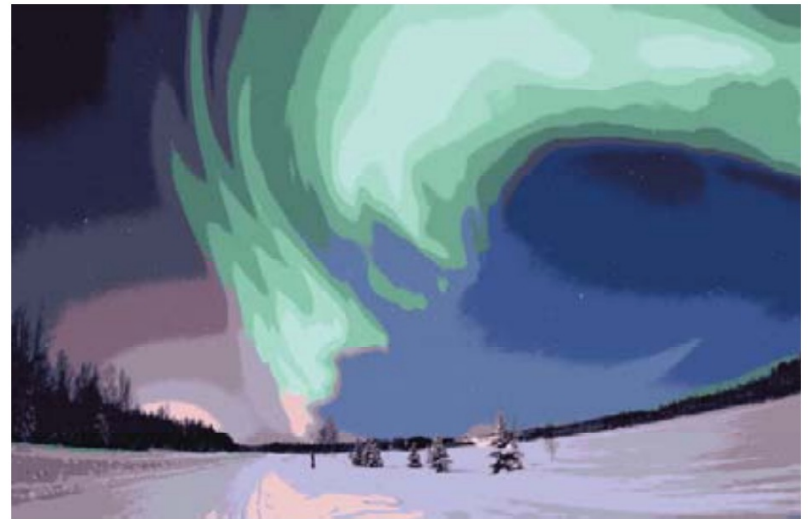


Cluster analysis and display of genome-wide expression patterns

[Michael B. Eisen](#),* [Paul T. Spellman](#),* [Patrick O. Brown](#),[†] and [David Botstein](#)^{*‡}

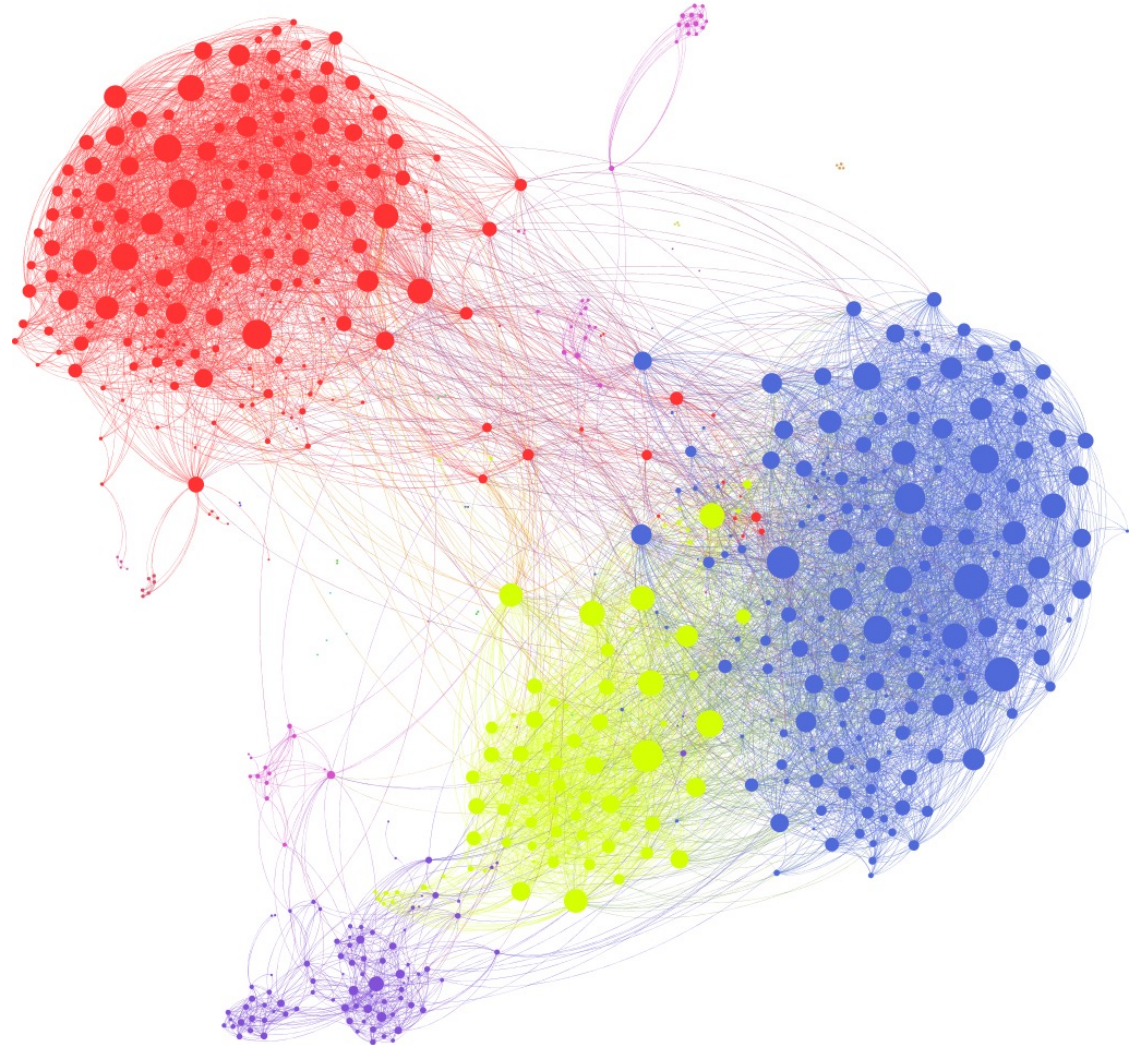
Applications of clustering

- Image segmentation: cluster similar regions of an image



Applications of clustering

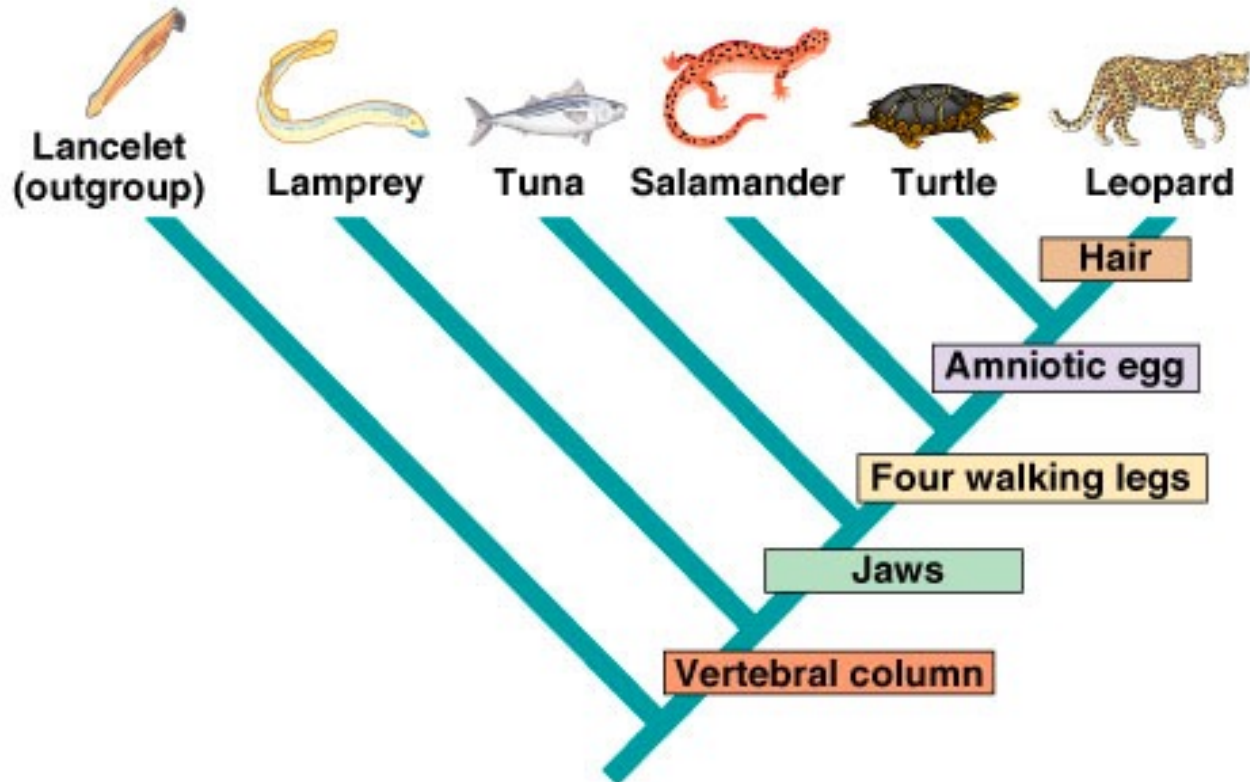
- Clustering in social graphs



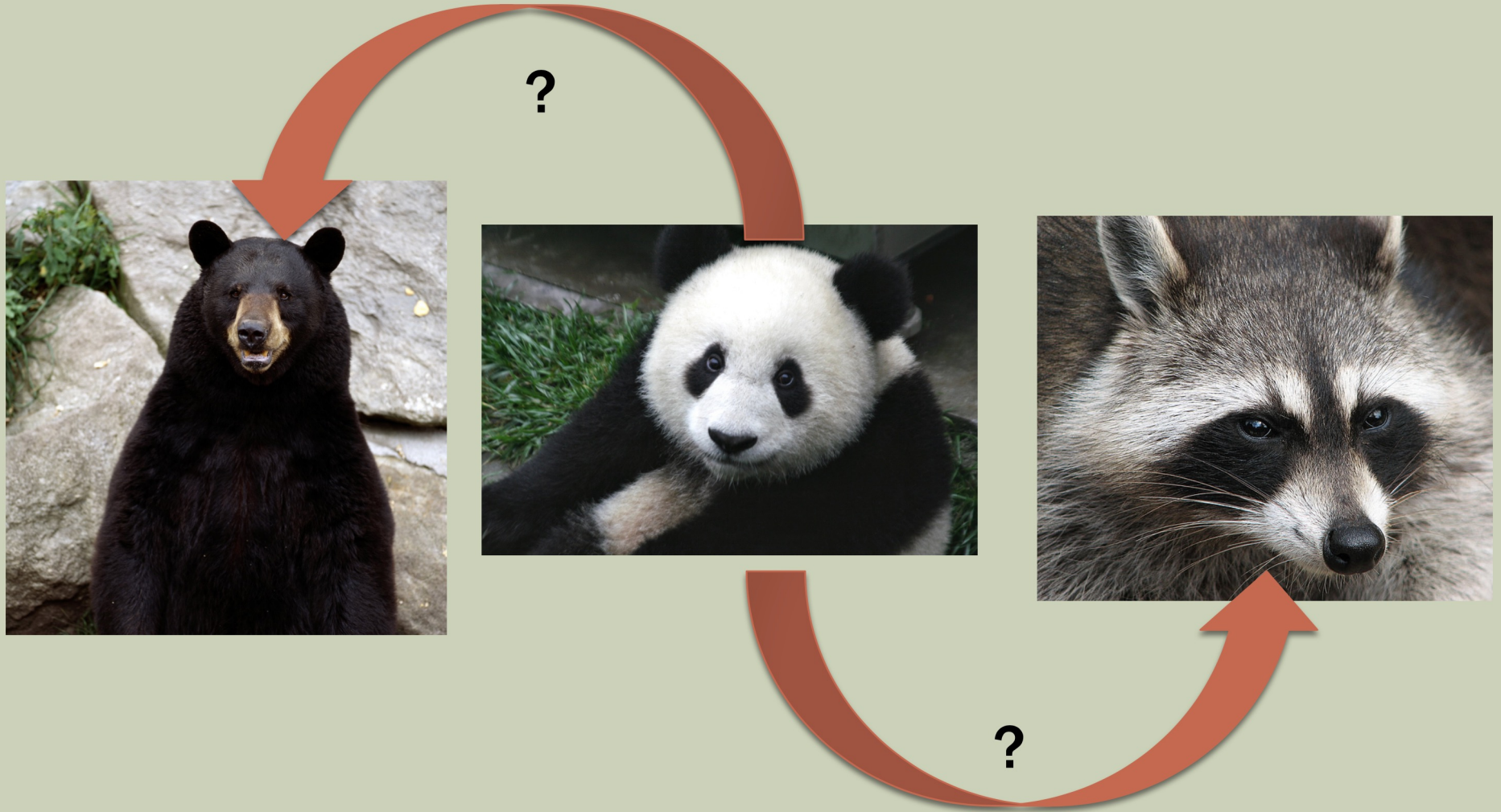
Two main types of clustering

- Flat/Partitional:
 - K-means
 - Gaussian mixture models
- Hierarchical:
 - Agglomerative: bottom-up
 - Divisive: top-down
 - Examples: UPGMA and Neighbor Joining

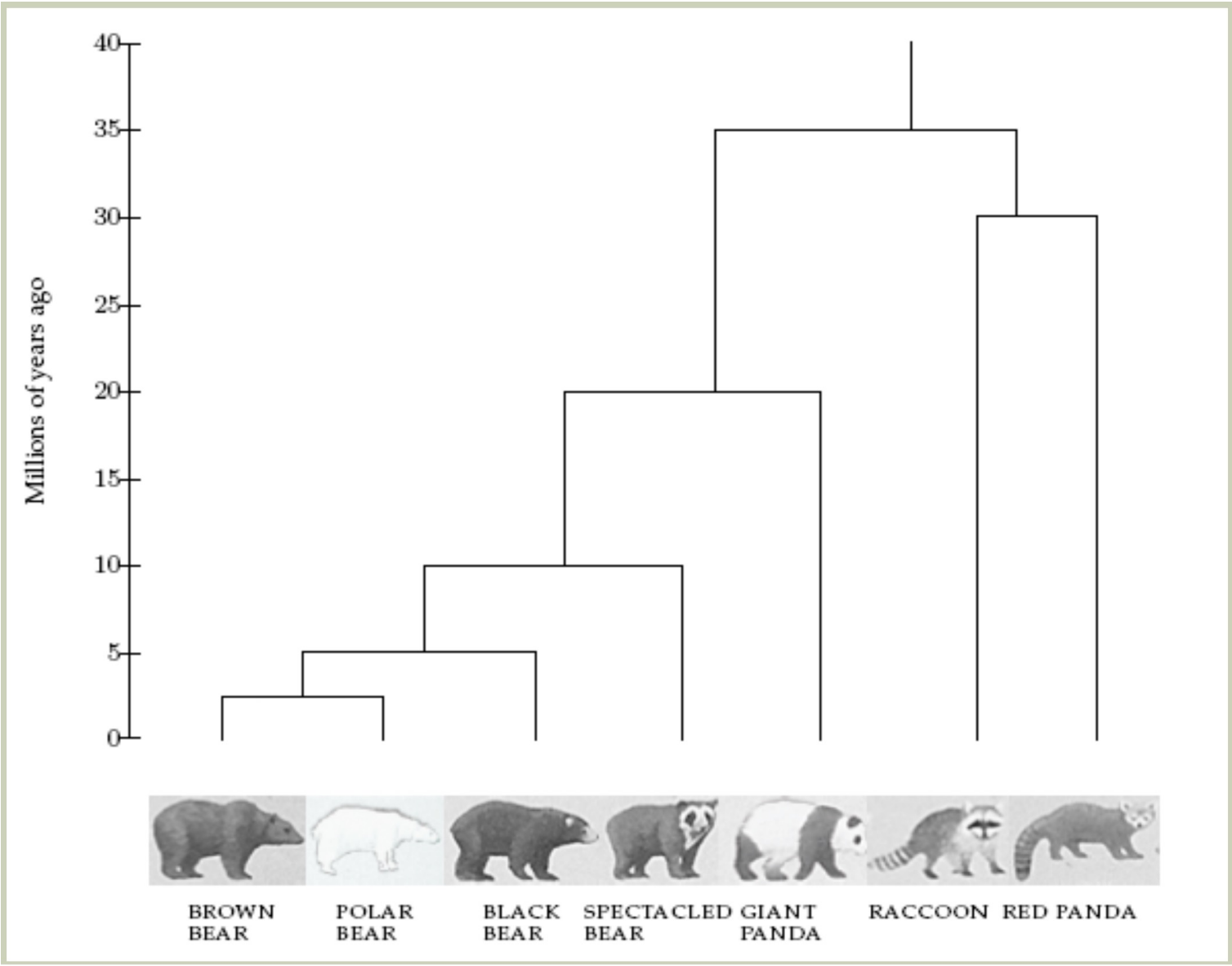
Hierarchical clustering example: trees



Are pandas more closely related to bears or raccoons?

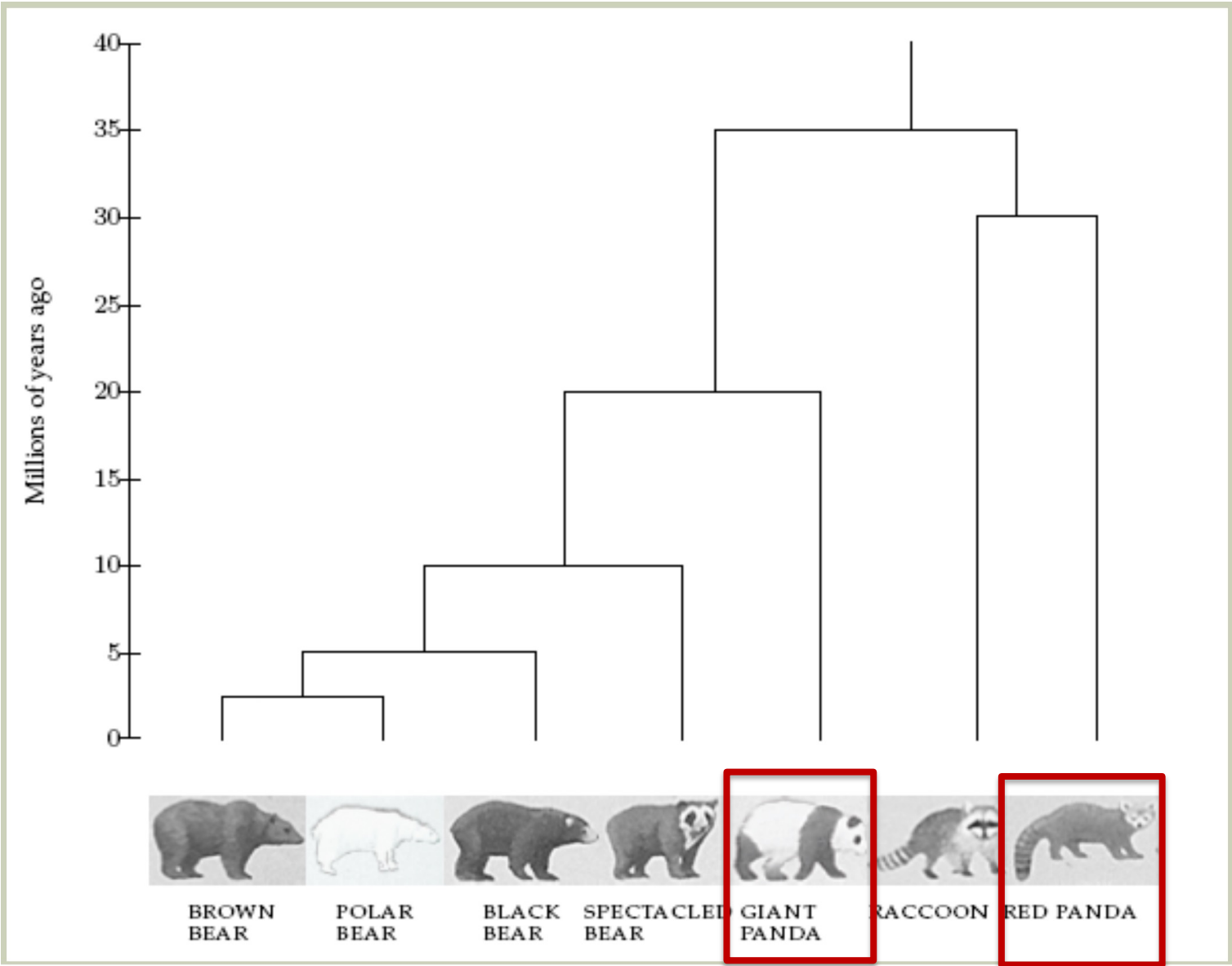


Are pandas more closely related to bears or raccoons?



Credit:
Ameet
Soni

Are pandas more closely related to bears or raccoons?



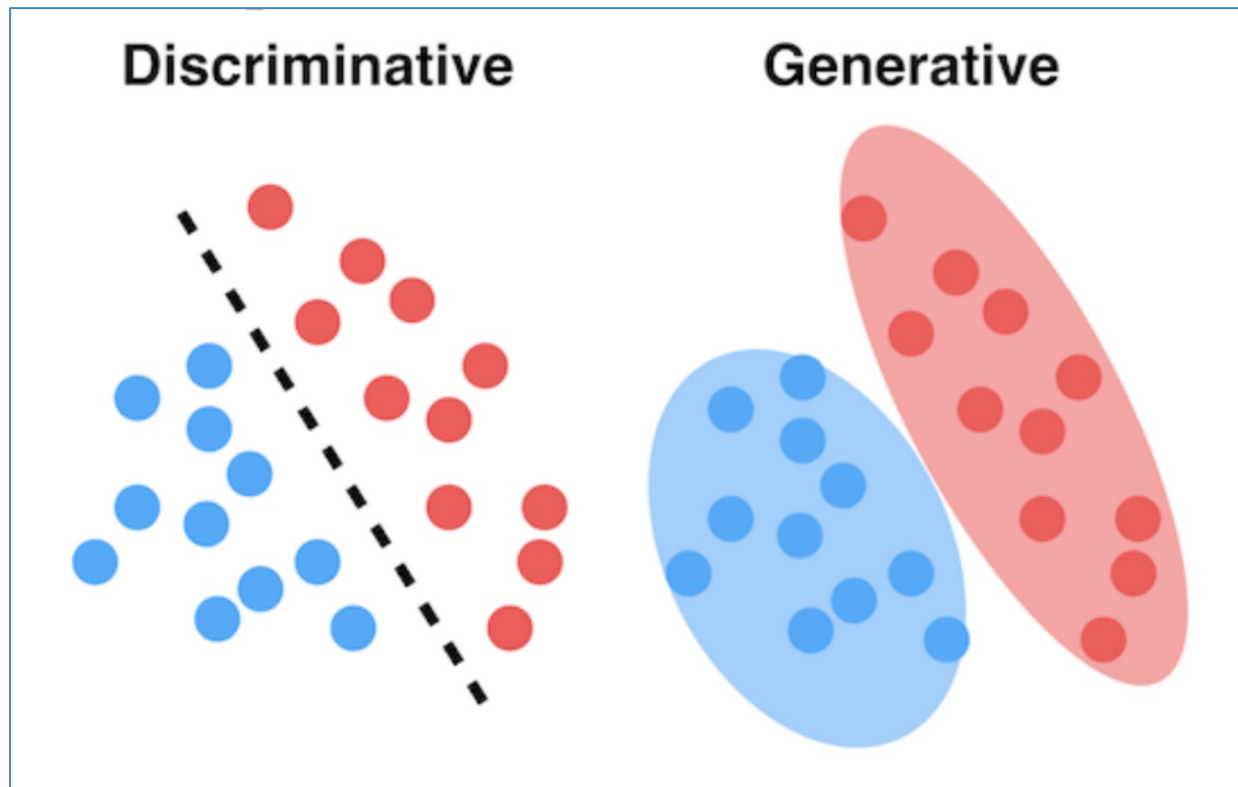
Credit:
Ameet
Soni

Outline for November 28

- Clustering overview
- **K-means**
- Gaussian Mixture Models (GMMs)

Discriminative vs. Generative

- Discriminative: finds a decision boundary
 - Logistic regression, K-means
- Generative: estimates probability distributions
 - Naïve Bayes, Gaussian Mixture Models



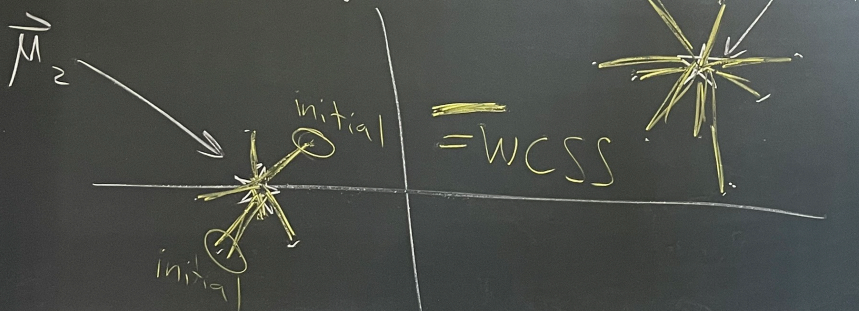
K-means

K clusters

minimize

$$WCSS = \sum_{k=1}^K \sum_{\vec{x}_i \in \mathcal{C}_k} \|\vec{x}_i - \vec{\mu}_k\|^2$$

within-cluster
Sum of squares



EM algorithm

E-step assignment

for \vec{x}_i , find closest cluster $\mathcal{C}_k^{(t)}$
mean $\Rightarrow \vec{x}_i \in \mathcal{C}_k$

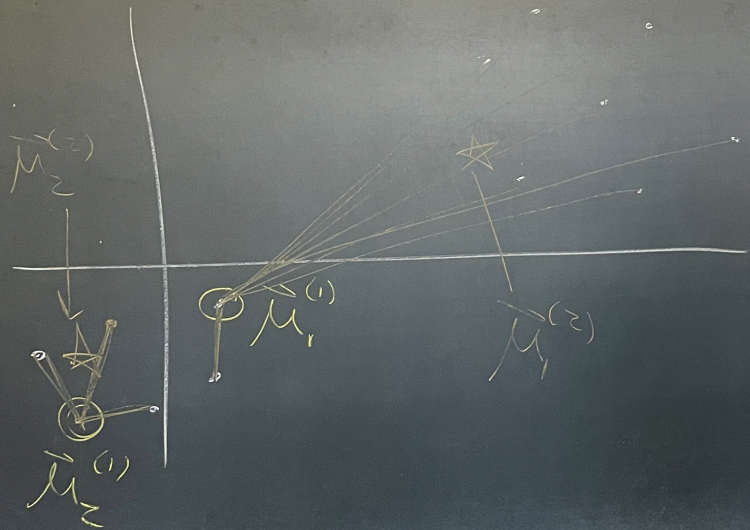
M-step parameter update

$$\mu_k^{(t+1)} = \frac{1}{|\mathcal{C}_k^{(t)}|} \sum_{\vec{x}_i \in \mathcal{C}_k^{(t)}} \vec{x}_i$$

Size of cluster k

initialization

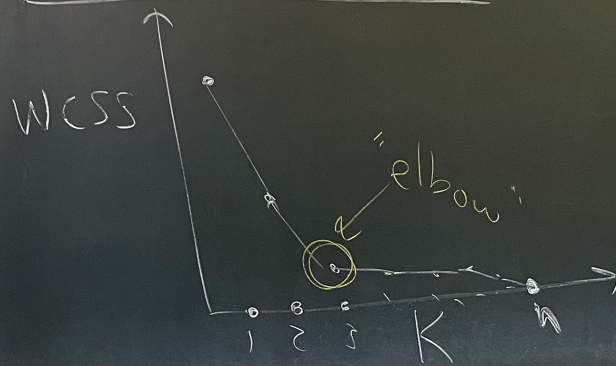
$\vec{\mu}$'s? choose randomly from our data points



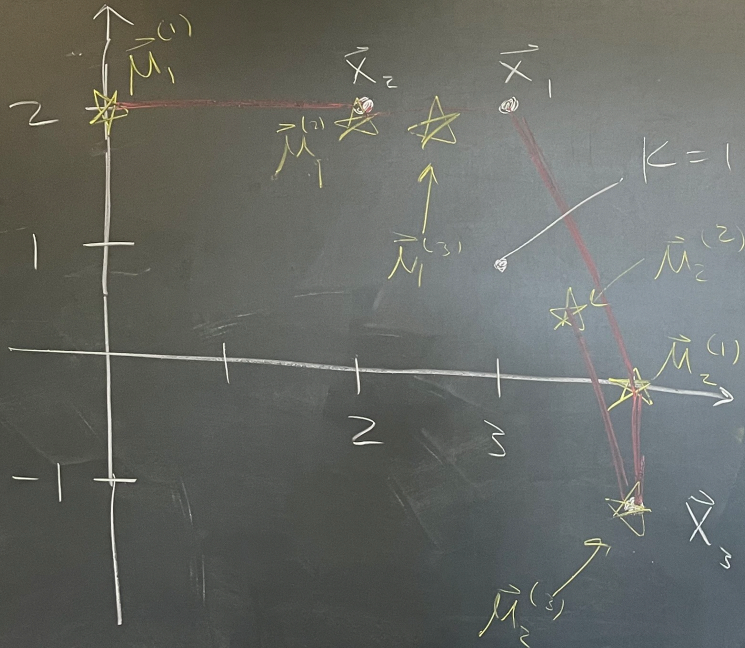
Stopping criteria

- no change in cluster membership
- max # iters exceeded (T)
- configuration you've see before.

how to choose K?



Handout 22



(a) **E-step**
 $C_1 = \{X_2\}$, $C_2 = \{X_1, X_3\}$

(b) **M-step**
 $\vec{\mu}_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$, $\vec{\mu}_2 = \begin{bmatrix} 3.5 \\ 0.5 \end{bmatrix}$

(c) **E-step**
 $C_1 = \{X_1, X_2\}$, $C_2 = \{X_3\}$

M-step
 $\mu_1^{(3)} = \begin{bmatrix} 2.5 \\ 2 \end{bmatrix}$, $\mu_2^{(3)} = \begin{bmatrix} 4 \\ -1 \end{bmatrix}$

(2) yes (monotonic)

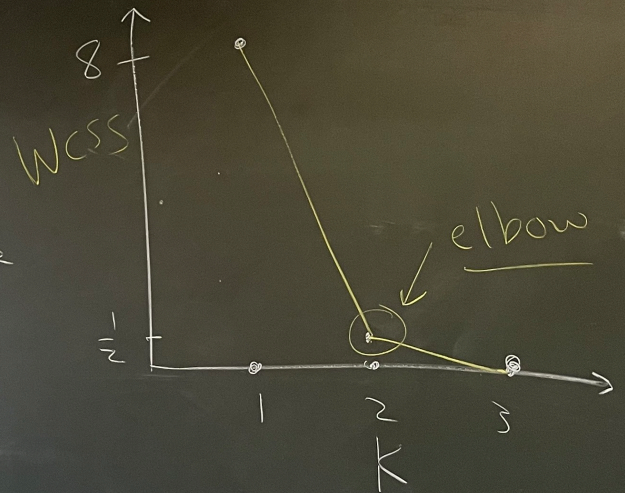
(3) $K=3$, $WCSS=0$

$$K=2, WCSS = \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 + 0^2 \\ = \frac{1}{2}$$

$$K=1, \bar{\mu}_1 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

$$X = \begin{bmatrix} 3 & 2 \\ 2 & 2 \\ 4 & 1 \end{bmatrix}$$

$$WCSS = 1^2 + (\sqrt{2})^2 + (\sqrt{5})^2 = 8$$



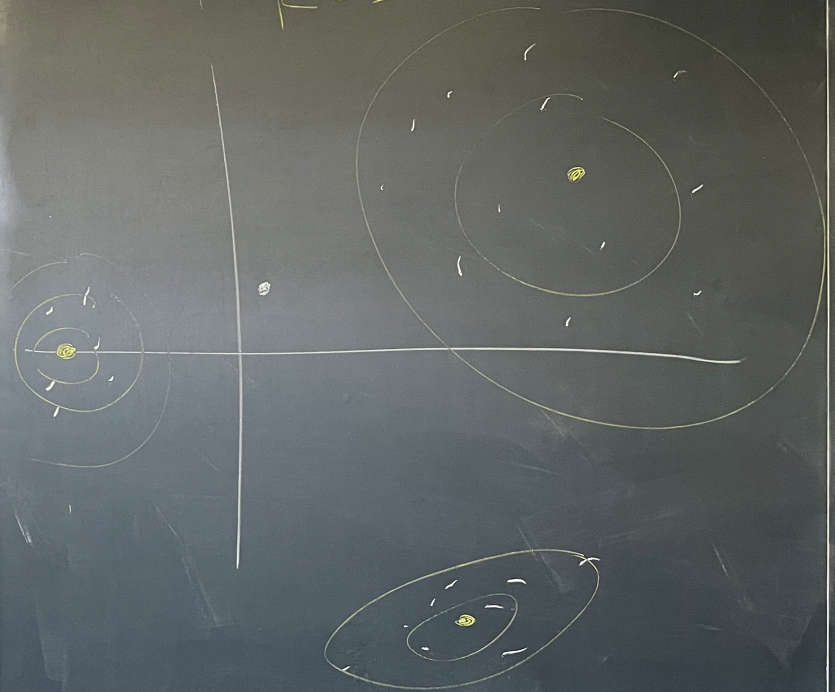
Outline for November 28

- Clustering overview
- K-means
- Gaussian Mixture Models (GMMs)

Problems with K-means

- * not generative (could not create a new data-point)
- * does not account for different cluster sizes & variances
- * does not allow points to belong to multiple clusters.

$K=3$



Gaussian Mixture Models

Likelihood \leftarrow maximize

$$p(\vec{x}) = \sum_{k=1}^K p(\vec{x}, z=k) =$$

cluster membership

BAYES

$$L(X) = \prod_{i=1}^n \sum_{k=1}^K \pi_k N(\vec{x}_i; \vec{\mu}_k, \sigma_k^2)$$

one pt \rightarrow \vec{x}_i

all data \rightarrow $\prod_{i=1}^n$

don't know \rightarrow $\sum_{k=1}^K$

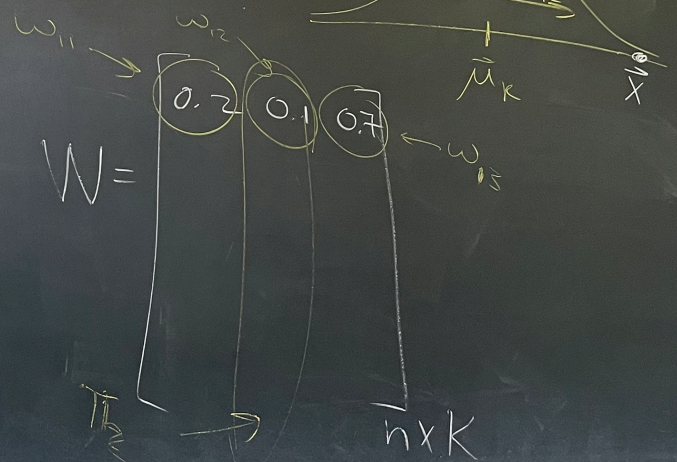
know \rightarrow $\vec{\mu}_k, \sigma_k^2$

(GMM)

prior over cluster sizes

$$\sum_{k=1}^K \pi_k p(\vec{x} | z=k)$$

Gaussian with $\vec{\mu}_k, \sigma_k^2$



EM for GMM

Initialization

• $\pi_k = \frac{1}{K}$ (uniform) (prior)

• $\vec{\mu}_k =$ choose random points

• $\sigma_k^2 =$ sample variance of all points closest to each mean

E-step 'Soft' assignment

$w_{ik} =$ Prob that \vec{x}_i came from cluster k

$$w_{ik} = p(k | \vec{x}_i) = \frac{p(k)p(\vec{x}_i | k)}{p(\vec{x}_i)}$$

$$= \frac{\pi_k \mathcal{N}(\vec{x}_i; \vec{\mu}_k, \sigma_k^2)}{\sum_{k'=1}^K \pi_{k'} \mathcal{N}(\vec{x}_i; \vec{\mu}_{k'}, \sigma_{k'}^2)}$$

Example of GMMs with different covariance constraints on the Iris flower data

