# CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024

HAVERFORD
COLLEGE

# Admin

- **Lab 8** due TODAY!
  - Extra office hours **4-5pm TODAY in H110** (Sara)

- **Midterm April 25** in class (next Thursday)
  - Can still do handout videos for extra credit! Up to 24 hours before exam

- **Project presentations**: last week of classes

- **Writeup** due by the end of finals period
  - May 11 for seniors (AND groups involving seniors)
  - May 17 for non-seniors

# Final Project Deliverables

- **Presentation**
  - Last week of classes
  - 12 min per pair
  - Peer feedback

- **Writeup**
  - In the style of a research paper
  - Text can be brief but should describe your motivation, hypotheses, data, methods, experiments, results, interpretation, and conclusions
  - At least 3 figures

README.md should have command lines for reproducibility!

# Lab 7 competition

- **Fejiro and Pranav**
  - Pooling layer, add FC layer with 1000 units, optimizer RMSprop
  - Test accuracy 66%

- **Gavin and Neha**
  - DenseNet layers
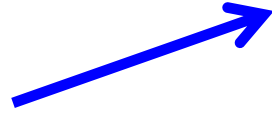  - Test accuracy 70%

# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

**Supervised Learning:** makes use of examples where we know the underlying "truth" (label/output)

**Unsupervised Learning:** Learn underlying structure or features without labeled training data

Note: **generative models** are typically unsupervised!

**Supervised learning** [hide]
(**classification** · **regression**)
Apprenticeship learning · Decision trees ·
Ensembles (Bagging · Boosting ·
Random forest) · $k$-NN · Linear regression ·
Naive Bayes · Artificial neural networks ·
Logistic regression · Perceptron ·
Relevance vector machine (RVM) ·
Support vector machine (SVM)

**Clustering** [hide]
BIRCH · CURE · Hierarchical · $k$-means ·
Fuzzy · Expectation–maximization (EM) ·
DBSCAN · OPTICS · Mean shift

**Dimensionality reduction** [hide]
Factor analysis · CCA · ICA · LDA · NMF · PCA
· PGD · t-SNE · SDL

**Structured prediction** [hide]
Graphical models (Bayes net ·
Conditional random field · Hidden Markov)

**Anomaly detection** [hide]
RANSAC · $k$-NN · Local outlier factor ·
Isolation forest

**Artificial neural network** [hide]
Autoencoder · Cognitive computing ·
Deep learning · DeepDream ·
Feedforward neural network ·
Recurrent neural network (LSTM · GRU · ESN
· reservoir computing) ·
Restricted Boltzmann machine · GAN ·
Diffusion model · SOM ·
Convolutional neural network (U-Net) ·
Transformer (Vision) · Mamba ·
Spiking neural network · Memtransistor ·
Electrochemical RAM (ECRAM)

**Reinforcement learning** [hide]
Q-learning · SARSA · Temporal difference (TD)
· Multi-agent (Self-play)

Image: wikipedia

# Unsupervised learning: main areas subject to debate

1) [Clustering]: group data points into clusters based on features only

2) [Dimensionality reduction]: remove feature correlation, compress data, visualize data

3) [Structured prediction]: model latent variables (example: Hidden Markov Models)

4) [Generative models]: learn latent structure in order to generate novel examples

# Applications of clustering
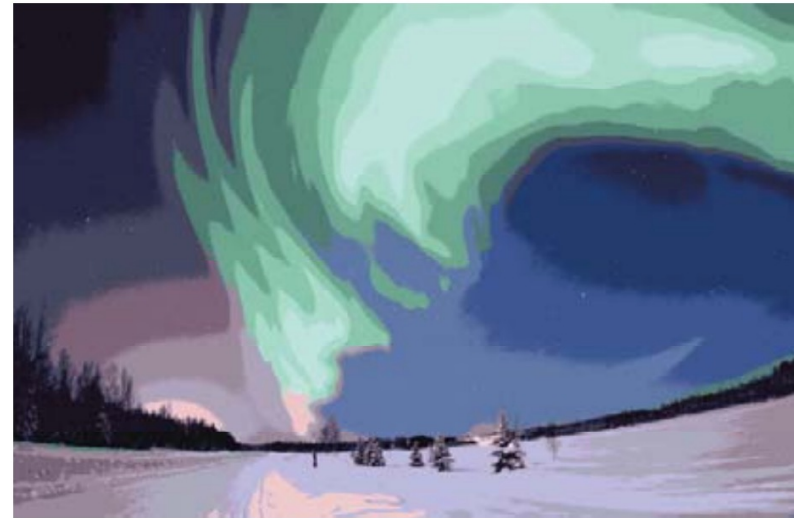
- Cluster genes with similar expression patterns

Cluster analysis and display of genome-wide expression patterns

Michael B. Eisen,[*] Paul T. Spellman,[*] Patrick O. Brown,[†] and David Botstein[*‡]
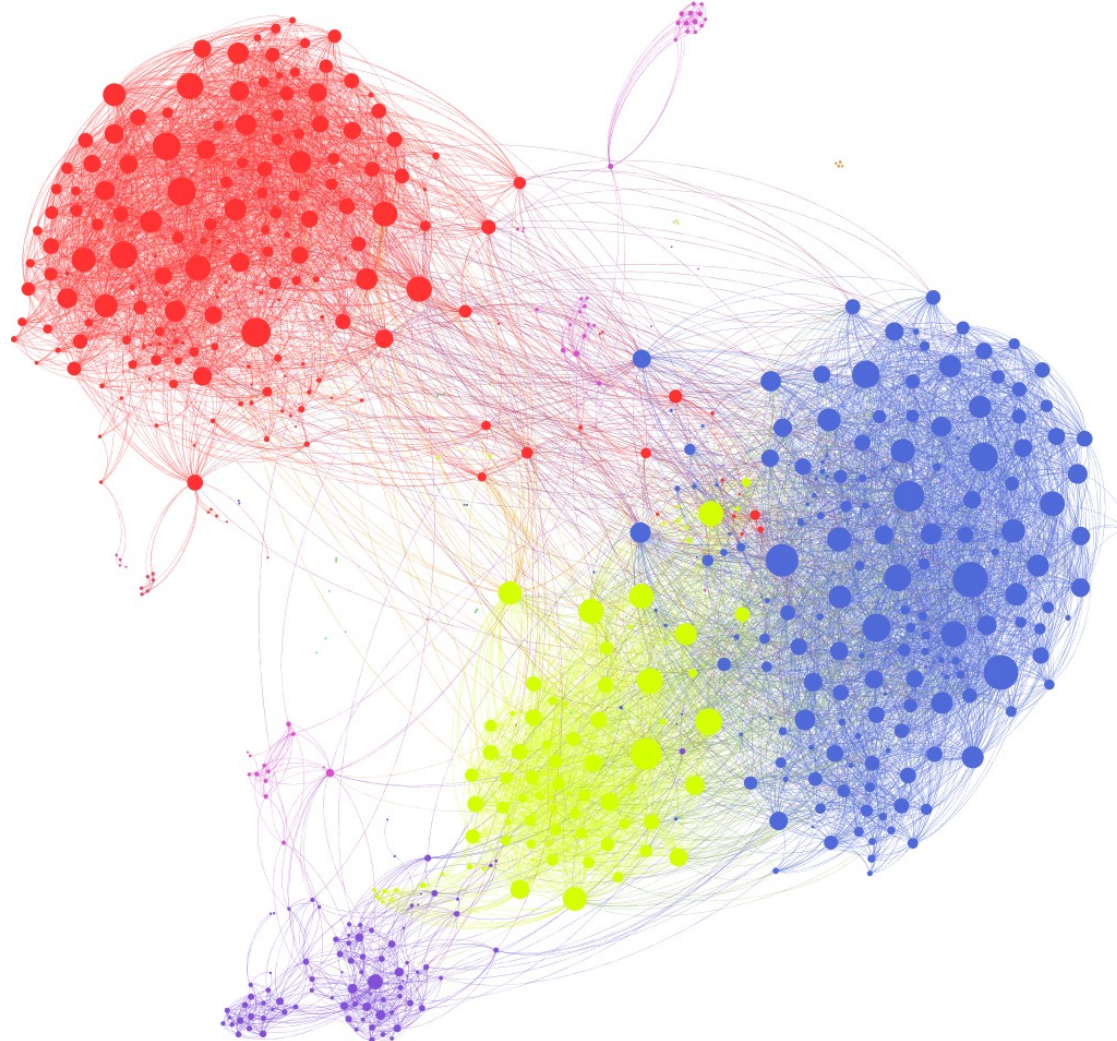
# Applications of clustering

- Image segmentation: cluster similar regions of an image

# Applications of clustering
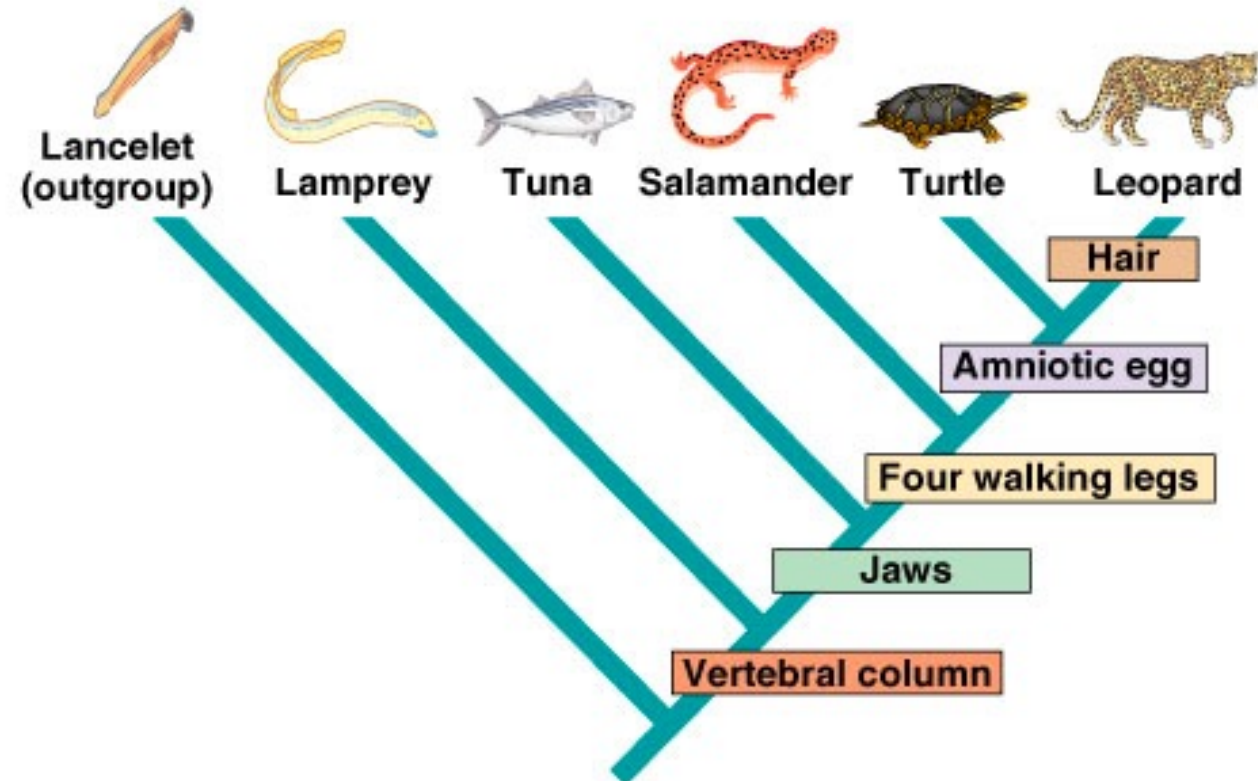


- Clustering in social graphs

# 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



Credit: Pearson Education, Benjamin Cummings

# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

# Review K-means from CS260

- Goal: learn about the structure in our data

- Goal: predict cluster membership for a new data point

- Method: minimize the within-cluster sum of squares (WCSS)

$$\mathcal{C} = \{\ \mathcal{C}_1, \quad \mathcal{C}_2 \quad \cdots \quad \mathcal{C}_K\ \}$$

$$\Downarrow \qquad \Downarrow \qquad\qquad \Downarrow$$

$$\vec{M}_1 \qquad \vec{M}_2 \qquad\qquad \vec{M}_K\ \} \text{ means}$$

$$\mathcal{C}_k = \{\ \vec{X}_7, \quad \vec{X}_{12}, \quad \vec{X}_{18}\ \}$$

Minimize

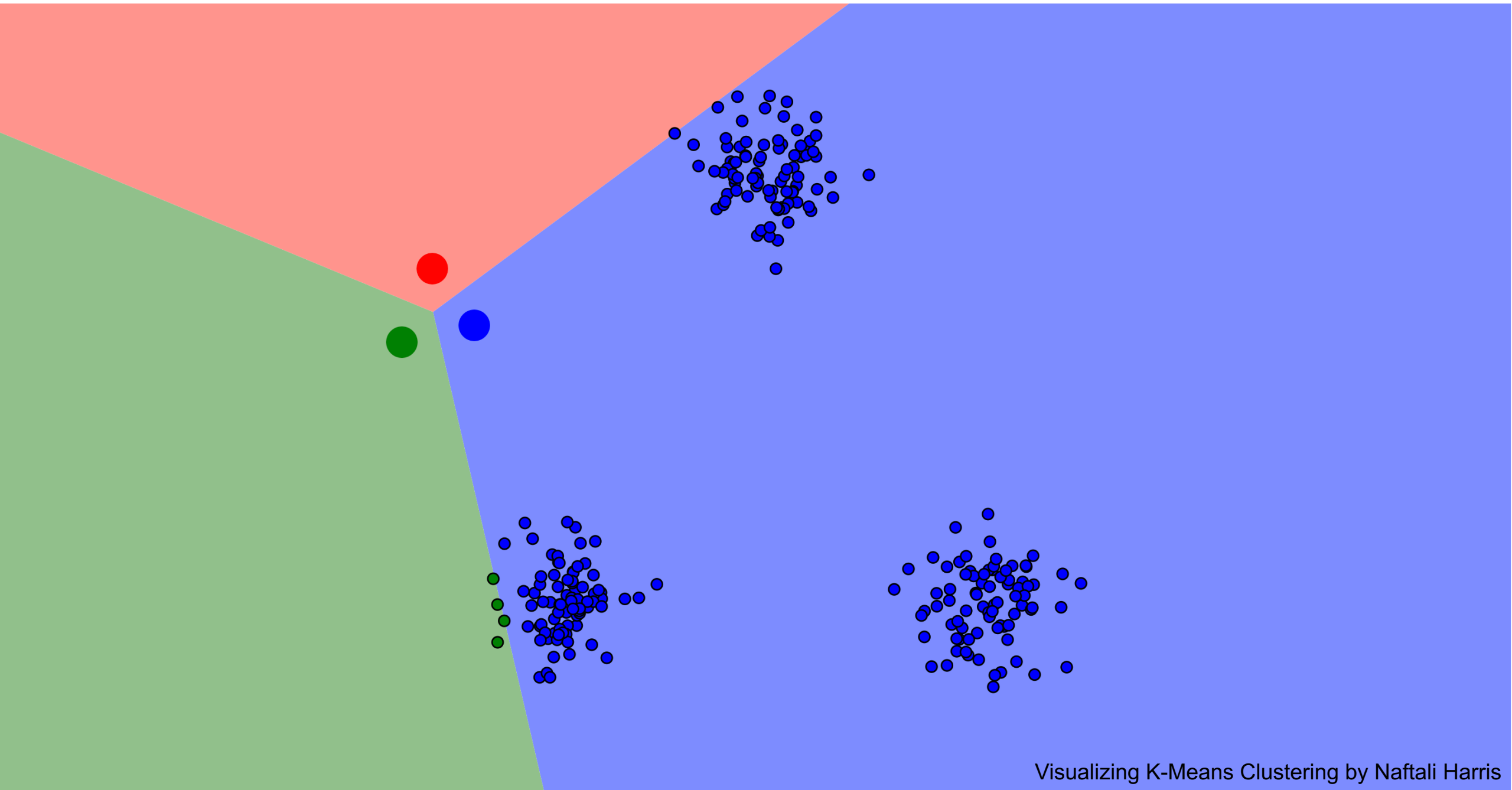$$\uparrow \text{WCSS} = \boxed{\sum_{k=1}^{K} \sum_{i \in \mathcal{C}_k} \|\vec{X}_i - \vec{M}_k\|^2}$$
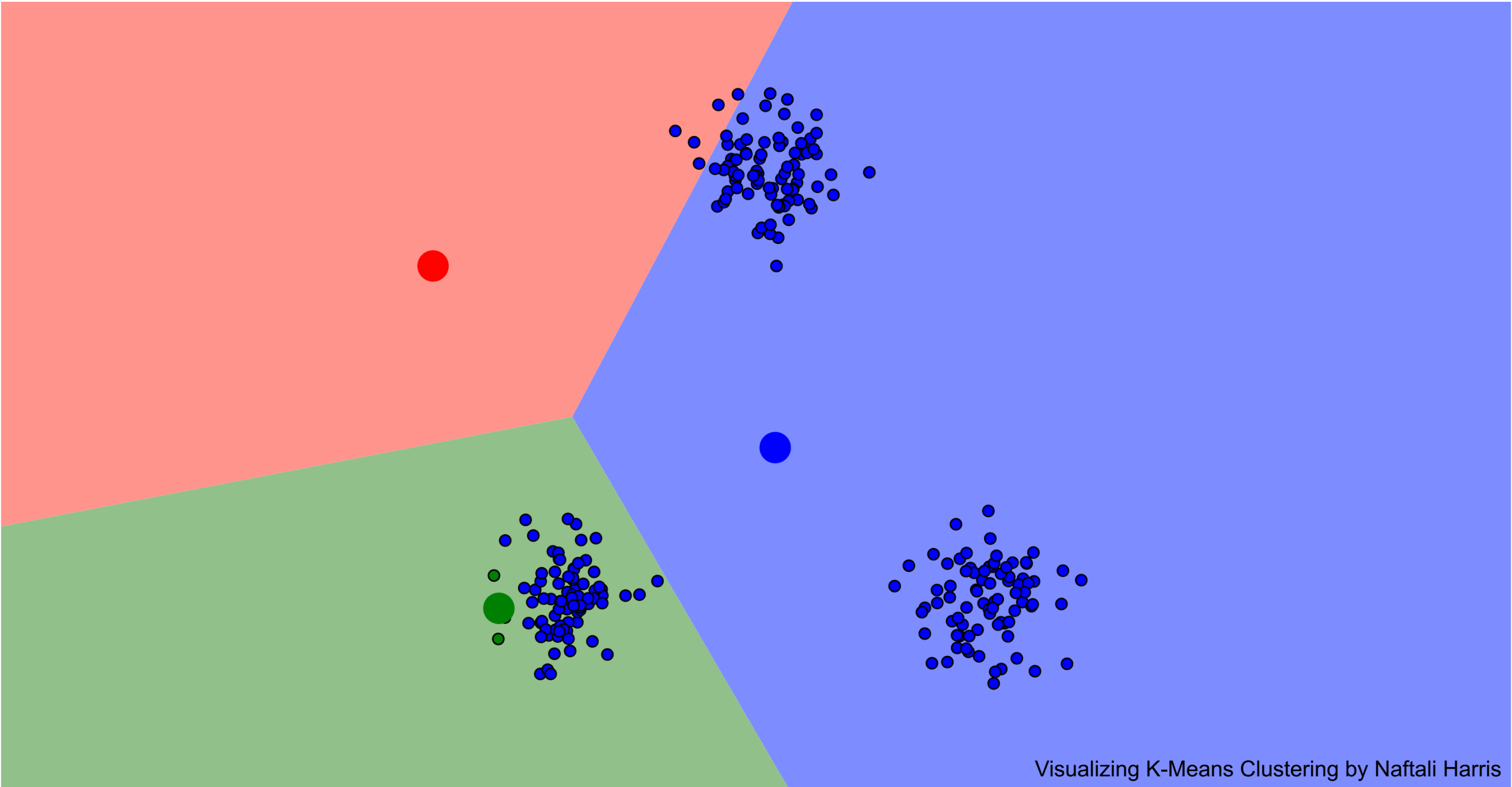
NP-hard

# Review K-means from CS260

- **Initialization**: choose K means (cluster centers) randomly

  – Usually from among the training data

- **E-step**: assign each point to the closest mean

- **M-step**: update the means as the cluster averages

Visualizing K-Means Clustering by Naftali Harris

Visualizing K-Means Clustering by Naftali Harris

Visualizing K-Means Clustering by Naftali Harris

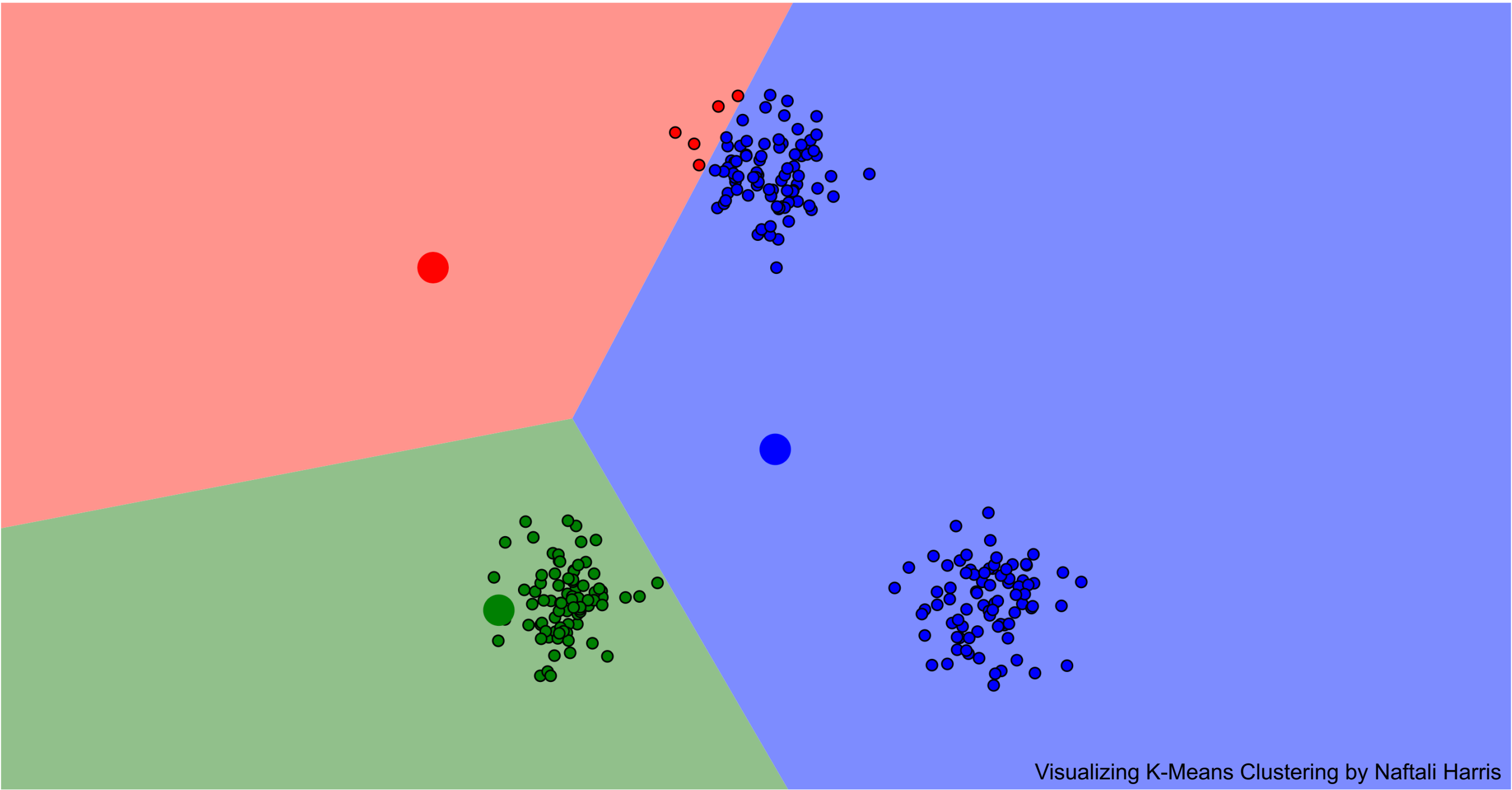Visualizing K-Means Clustering by Naftali Harris

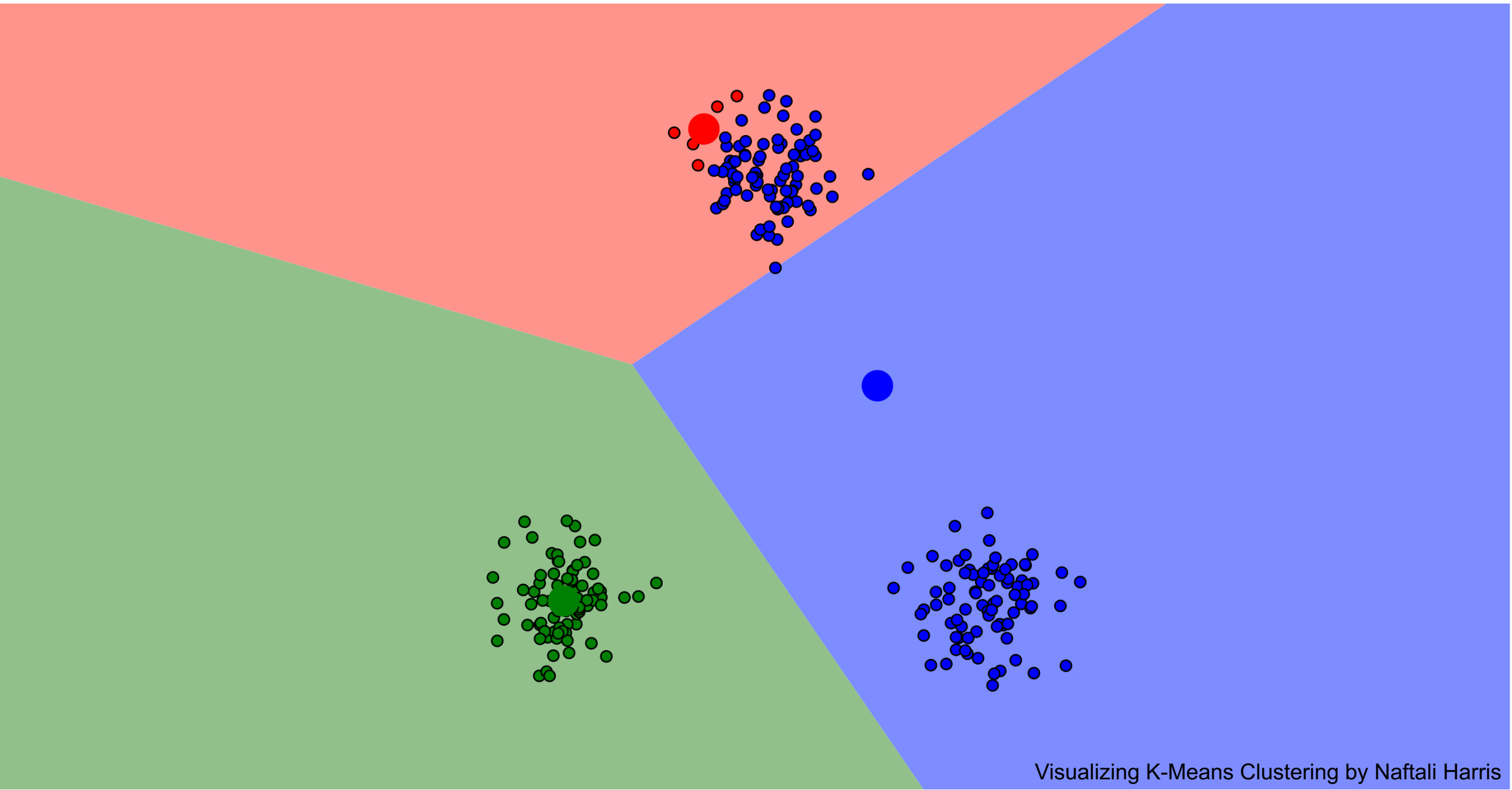Visualizing K-Means Clustering by Naftali Harris

Visualizing K-Means Clustering by Naftali Harris

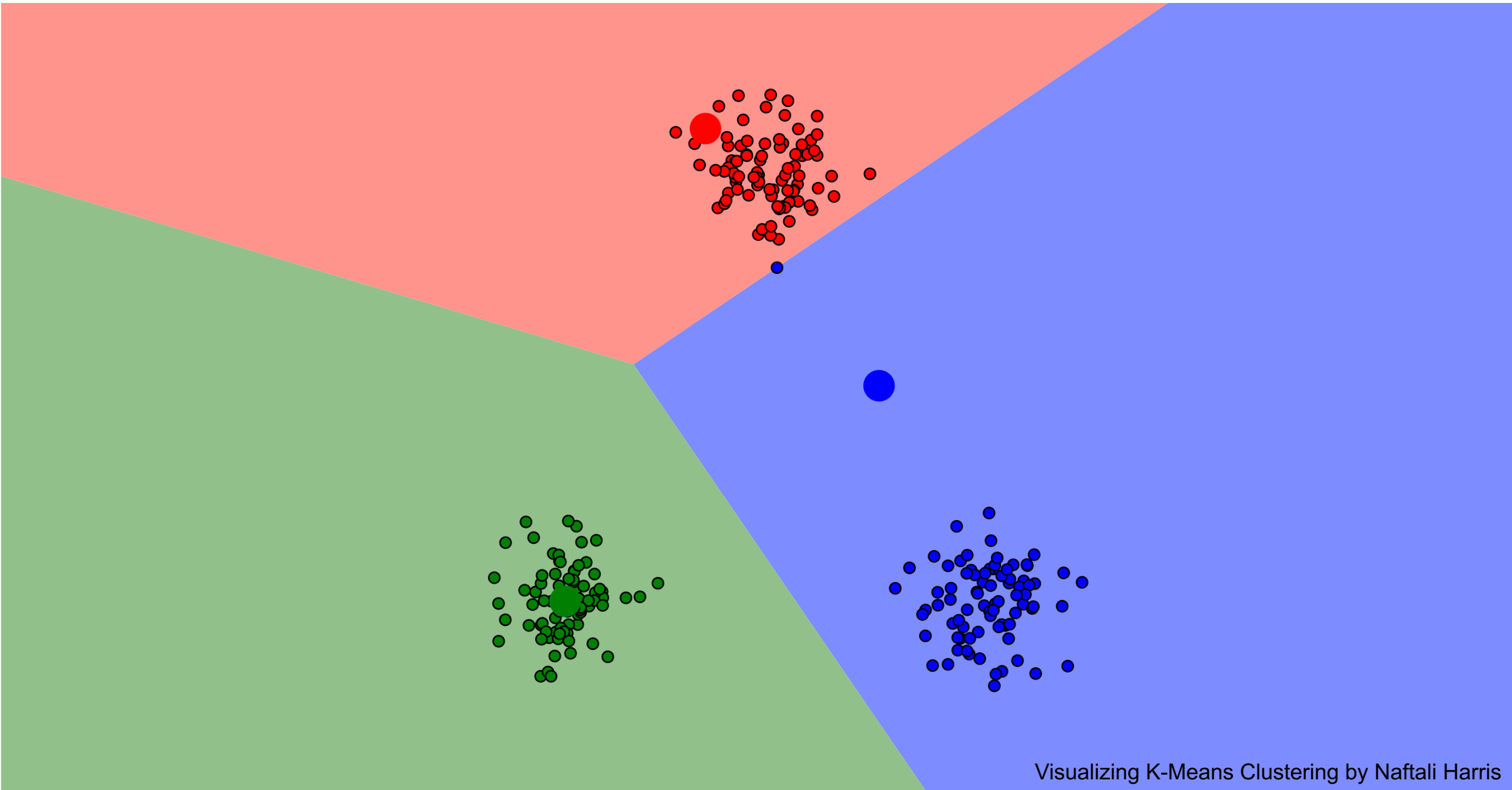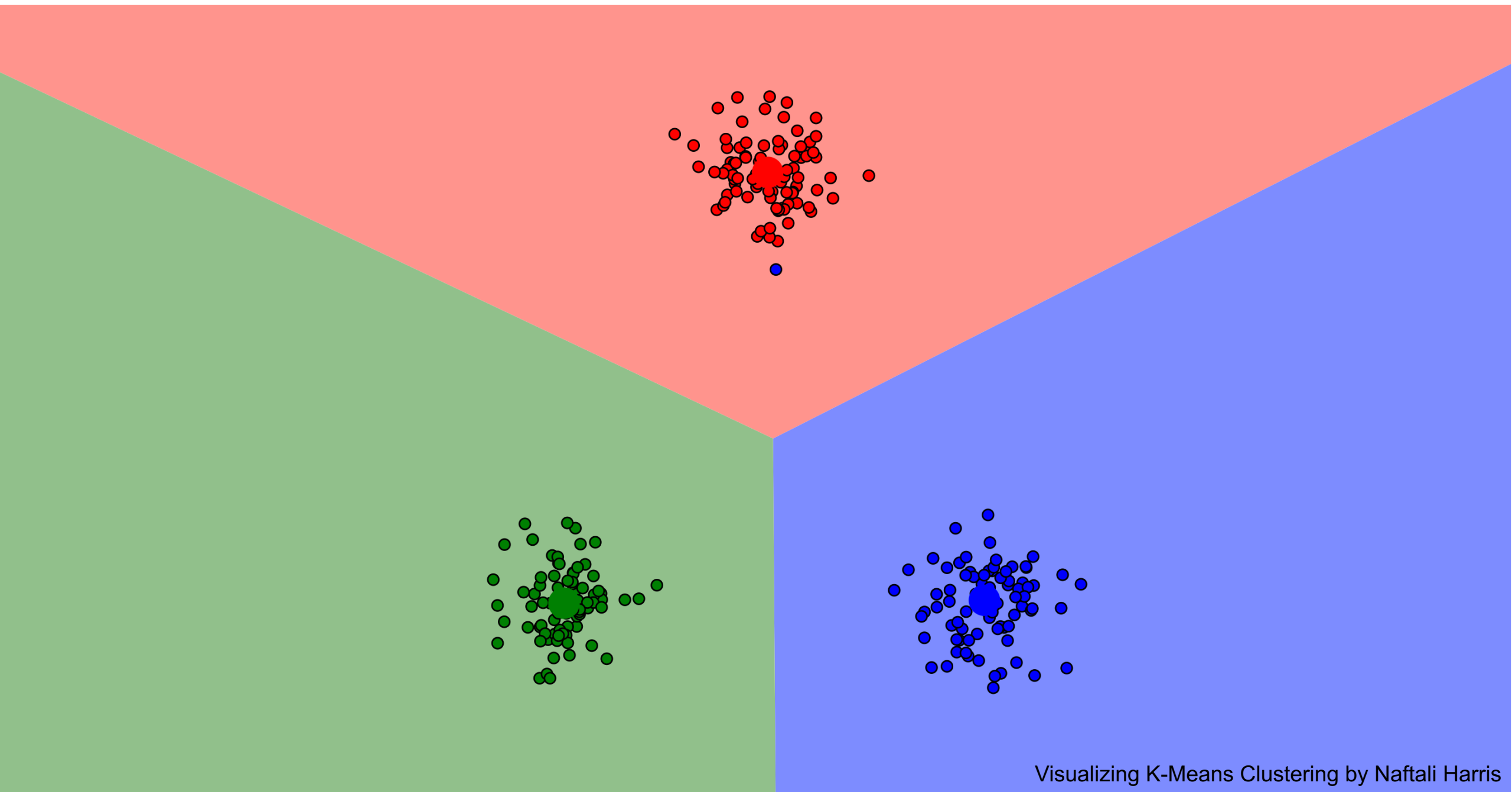Visualizing K-Means Clustering by Naftali Harris

Visualizing K-Means Clustering by Naftali Harris

# How to choose the number of clusters K?

- A larger K will always reduce the within-cluster sum of squares (WCSS)!

- Try a variety of K and choose the K value where the decreases in WCSS plateau

- **Question**: what is the WCSS when K=n? (the number of training points)



Image: Datanovia

# K-means stopping criteria

- No cluster membership changes

- Max number of iterations exceeded

- See a configuration you've seen before (cycle)

# Discriminative vs. Generative

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



Figure: Ameet Soni

# Problems with K-means

- Not generative (could not create a new data point)

- Does not account for different cluster sizes and variances

- Does not allow points to belong to multiple clusters

# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

# Gaussian Mixture Models GMM

$$\text{(K clusters)}$$

## Likelihood

$$p(\vec{x}) = \sum_{k=1}^{K} p(\vec{x}, \underset{\uparrow}{z}=k) = \sum_{k=1}^{K} \underbrace{p(z=k)}_{\pi_k} \underbrace{p(\vec{x}|z=k)}_{\substack{\text{given cluster,} \\ \text{prob of } \vec{x}}}$$

cluster membership

$\pi_k$ — size of cluster k

$\longrightarrow$ **BAYES**

$$L(X) = \prod_{i=1}^{n} \sum_{k=1}^{K} \pi_k \underbrace{\mathcal{N}(\vec{x}_i ; \underset{\substack{\uparrow \\ \text{mean}}}{\vec{\mu}_k}, \underset{\text{variance}}{\sigma_k^2})}_{\text{normal distribution}}$$

all train data

## Goal:

find

$$\boxed{\pi_k, \ \vec{\mu}_k, \ \sigma_k^2}$$

that maximize likelihood!

# EM algorithm

- $\pi_k = $ prob of cluster $k$
  $$= \frac{1}{k} \text{ to start}$$

- $\vec{M}_k = $ mean of cluster $k$
  $$= \text{random data pt to start}$$

- $\sigma_k^2 = $ variance of cluster $k$

  $$= \text{Sample variance of all points}$$
  $$\text{closest to } \vec{M}_k$$

let

$W_{ik}$

$\begin{bmatrix} W_{11} & W_{12} & W_{13} \end{bmatrix}$

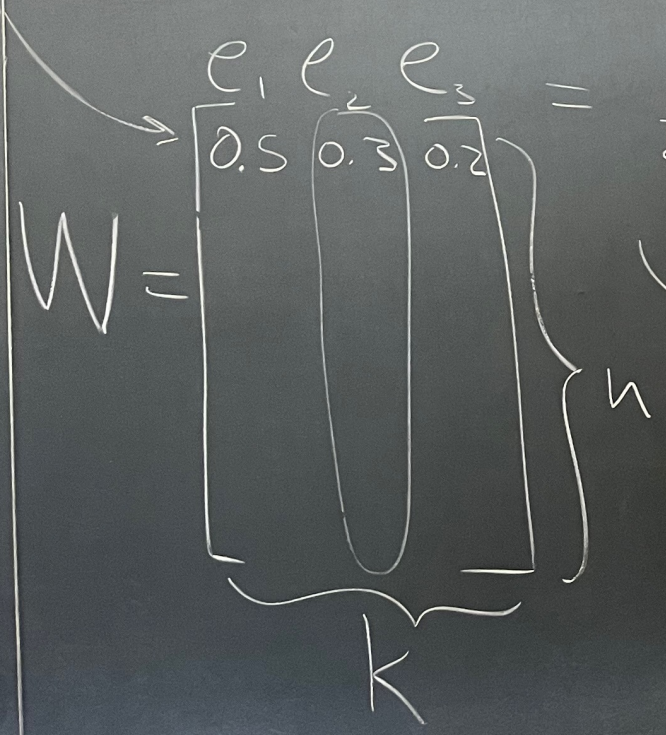$W = \begin{bmatrix} 0. \\ \end{bmatrix}$

$$\boxed{\text{E-step}} \text{ (soft assignment)}$$
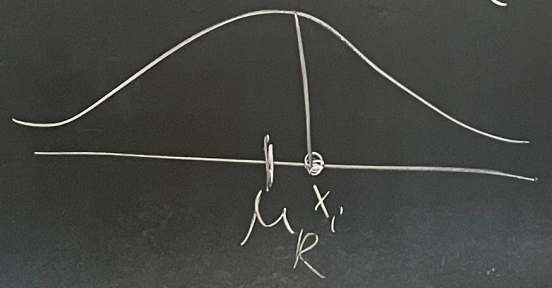
let $W_{ik} = $ prob that $\vec{x}_i$ came from cluster $k$

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

$W_{13}$

$$\begin{array}{ccc} e_1 & e_2 & e_3 \end{array} = \frac{\pi_k\, N(\vec{x}_i\,;\,\vec{M}_k,\,\sigma_k^2)}{\underbrace{\sum\limits_{k'} \pi_{k'}\, N(\vec{x}_i\,;\,\vec{M}_{k'},\,\sigma_{k'}^2)}_{\text{normalize}}}$$

$$W = \left[\begin{array}{|c|c|c|} \hline 0.5 & 0.3 & 0.2 \\ \hline & & \\ & & \\ & & \\ \hline \end{array}\right\} n$$

$$\underbrace{\qquad\qquad}_{K}$$

$\mu_k \quad x_i$

$$\boxed{\text{M-step}} \quad \text{let} \quad M_k = \sum\limits_{i=1}^{n} W_{ik}$$

$$\boxed{\pi_k = \frac{M_k}{n}}$$

$\underbrace{\qquad}$ "#" of points in cluster $k$

$$\boxed{\vec{M}_k = \frac{1}{M_k} \sum\limits_{i=1}^{n} W_{ik}\,\vec{x}_i}$$

$$\sigma_k^2 = \text{weighted sample variance}$$

$$\begin{array}{c} \stackrel{\wedge}{\approx} \\ \\ 0 \end{array} \begin{bmatrix} 0.01 & 0.49 & 0.5 \\ \\ 0.01 & 0.99 & 0.0 \end{bmatrix}$$

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

*Density Estimation with Gaussian Mixture Models*



(a) Dataset.

(b) Negative log-likelihood.

(c) EM initialization.

(d) EM after one iteration.

(e) EM after 10 iterations.

(f) EM after 62 iterations.

Figure 11.7 from MML textbook

# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

# Principal Components Analysis (PCA)

- Transforms *p*-dimensional data so that the new first dimension explains as much of the variation as possible, the new second explains as much of the remaining variation as possible, and so on

- Typically, we look at the first few dimensions of the transformed data and use as a means of dimensionality reduction and visualization

- PCA is a linear transformation

- PCA is often used for:
  - Data visualization
  - Infer qualitative relationships between groups

# Principal component analysis

# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

# Traditional Autoencoder

- Project data in a lower dimension (encoder), then reconstruct using a decoder
- Loss function: minimize reconstruction error

Traditional Autoencoder on MNIST

Optimizing purely for reconstruction loss

Intuitively Understanding Variational Autoencoders by Irhum Shafkat

# Variational Autoencoder: latent space has a set of means and variances



Variational Autoencoder

# Example of a trained VAE

Output
$\mu$

$[0.1, 1.2, 0.2, 0.8, \ldots]$

Output
$\sigma$

$[0.2, 0.5, 0.8, 1.3, \ldots]$

Intermediate
X

$[X_1 \sim N(0.1,\ 0.2^2),\ X_2 \sim N(1.2,\ 0.5^2),\ X_3 \sim N(0.2,\ 0.8^2),\ X_4 \sim N(0.8,\ 1.3^2), \ldots]$

sample

Sampled
vector

$[0.28, 1.65, 0.92, 1.98, \ldots]$

Stochastically generating encoding vectors

# VAE: don't want "spaces" between our clusters



What we require

What we may inadvertently end up with

Optimizing using pure KL divergence loss

Optimizing KL divergence only

Optimizing using both reconstruction loss and KL divergence loss

# Full VAE loss

Intuitively Understanding Variational Autoencoders by Irhum Shafkat

# Biology Example: popvae



**Figure 1** A schematic of the VAE architecture.

# Biology Example: popvae



**Figure 2** PCA axes 1–8 (left) and popvae run at default settings (right) for 100,000 random SNPs from chromosome 1

*Figure 17-5. Fashion MNIST visualization using an autoencoder followed by t-SNE*
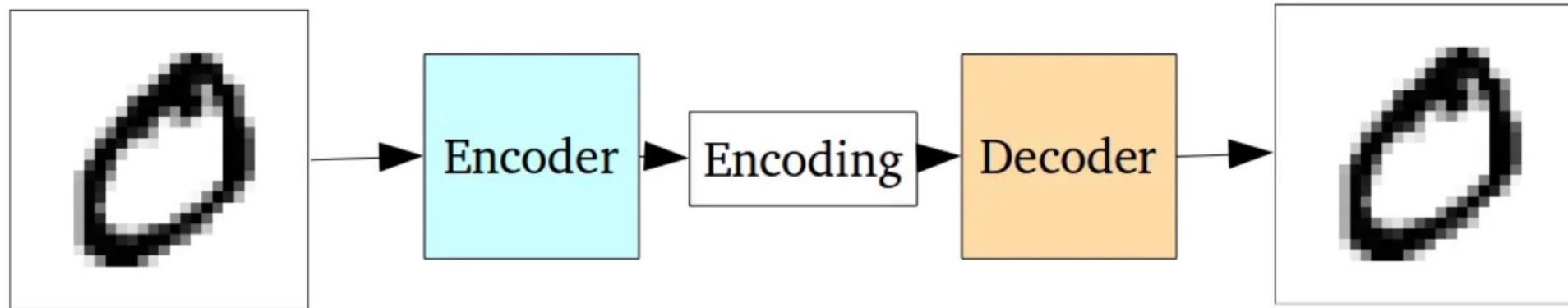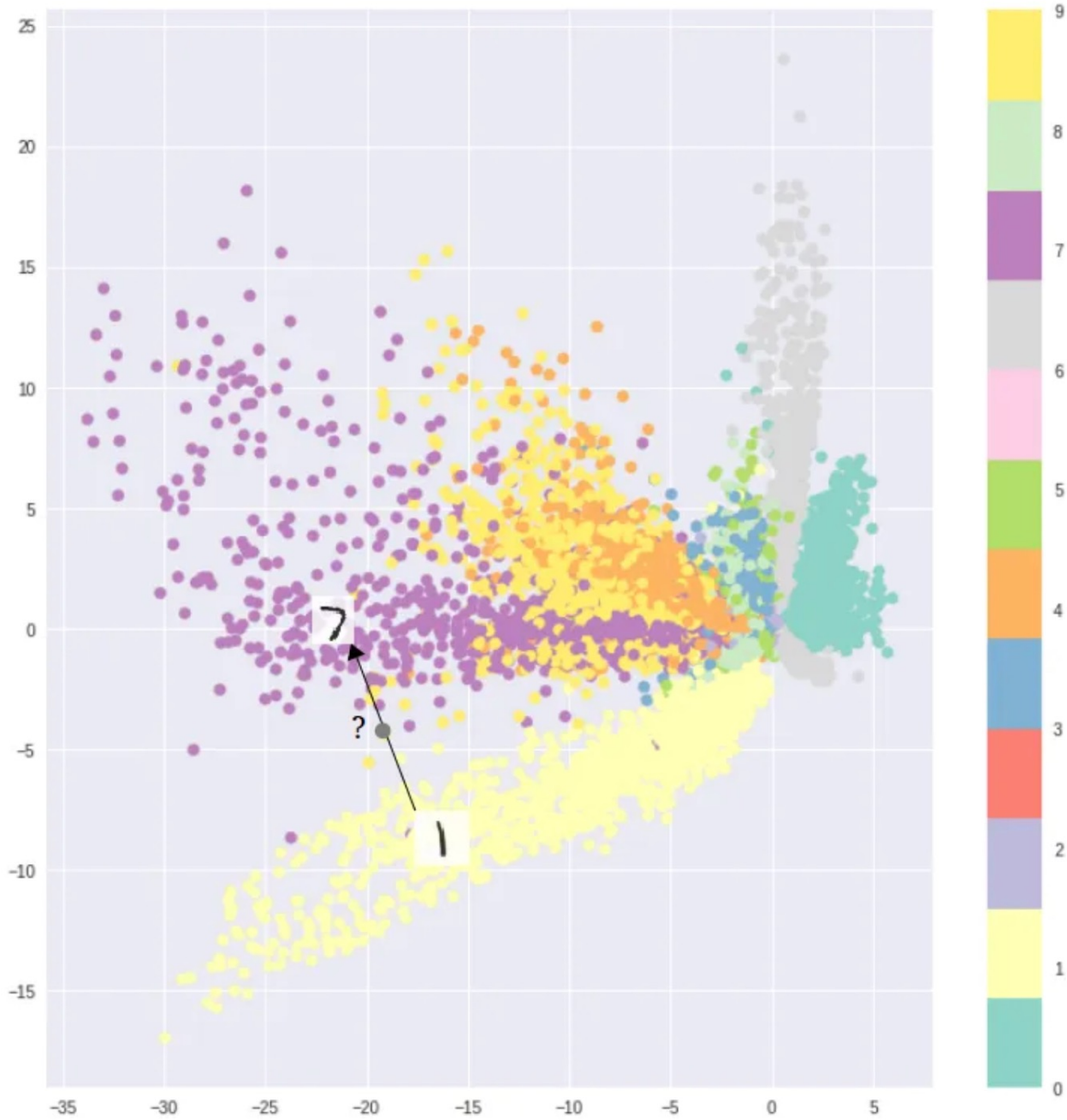
# Outline for April 18

- Introduction to unsupervised learning

- Review K-means (from CS260)

- Gaussian Mixture Models (GMMs)

- Review PCA (from CS260)

- Autoencoders

- Variational Autoencoders (VAEs)

- Hierarchical clustering (if time)

# Are pandas more closely related to bears or raccoons?

# UPGMA and Neighbor Joining

- Start with a dissimilarity map between examples (symmetric matrix)

- Say our examples are: A,B,C,D,E

| $\delta$ | A | B | C | D | E |
|---|---|---|---|---|---|
| A | 0 | 1 | 3 | 6 | 6 |
| B |   | 0 | 2 | 5 | 5 |
| C |   |   | 0 | 5 | 5 |
| D |   |   |   | 0 | 2 |
| E |   |   |   |   | 0 |

# Hierarchical clustering example (UPGMA)

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

A  D  B  F  G  C  E

Figure: Dr. Richard Edwards

# Hierarchical clustering example (UPGMA)

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

|    | A | BF | C | D | E |
|----|---|----|---|---|---|
| BF | 18.50 |   |   |   |   |
| C  | 27.00 | 31.50 |   |   |   |
| D  | 8.00 | 17.50 | 26.00 |   |   |
| E  | 33.00 | 35.50 | 41.00 | 31.00 |   |
| G  | 13.00 | 12.50 | 29.00 | 14.00 | 28.00 |

A  D  B  F  G  C  E

0.5      0.5

0.0
0.5

0.5

# Hierarchical clustering example (UPGMA)

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

|    | A | BF | C | D | E |
|----|---|-----|---|---|---|
| BF | 18.50 |   |   |   |   |
| C  | 27.00 | 31.50 |   |   |   |
| D  | 8.00 | 17.50 | 26.00 |   |   |
| E  | 33.00 | 35.50 | 41.00 | 31.00 |   |
| G  | 13.00 | 12.50 | 29.00 | 14.00 | 28.00 |

|    | AD | BF | C | E |
|----|-----|-----|---|---|
| BF | 18.00 |   |   |   |
| C  | 26.50 | 31.50 |   |   |
| E  | 32.00 | 35.50 | 41.00 |   |
| G  | 13.50 | 12.50 | 29.00 | 28.00 |

Figure: Dr. Richard Edwards

# Hierarchical clustering example (UPGMA)

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

|   | A | BF | C | D | E |
|---|---|---|---|---|---|
| BF | 18.50 |   |   |   |   |
| C | 27.00 | 31.50 |   |   |   |
| D | 8.00 | 17.50 | 26.00 |   |   |
| E | 33.00 | 35.50 | 41.00 | 31.00 |   |
| G | 13.00 | 12.50 | 29.00 | 14.00 | 28.00 |

|   | AD | BF | C | E |
|---|---|---|---|---|
| BF | 18.00 |   |   |   |
| C | 26.50 | 31.50 |   |   |
| E | 32.00 | 35.50 | 41.00 |   |
| G | 13.50 | 12.50 | 29.00 | 28.00 |

|   | AD | BFG | C |
|---|---|---|---|
| BFG | 16.50 |   |   |
| C | 26.50 | 30.67 |   |
| E | 32.00 | 33.00 | 41.00 |

Figure: Dr. Richard Edwards

# Hierarchical clustering example (UPGMA)



|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

|   | A | BF | C | D | E |
|---|---|---|---|---|---|
| BF | 18.50 |   |   |   |   |
| C | 27.00 | 31.50 |   |   |   |
| D | 8.00 | 17.50 | 26.00 |   |   |
| E | 33.00 | 35.50 | 41.00 | 31.00 |   |
| G | 13.00 | 12.50 | 29.00 | 14.00 | 28.00 |

|   | AD | BF | C | E |
|---|---|---|---|---|
| BF | 18.00 |   |   |   |
| C | 26.50 | 31.50 |   |   |
| E | 32.00 | 35.50 | 41.00 |   |
| G | 13.50 | 12.50 | 29.00 | 28.00 |

|   | AD | BFG | C |
|---|---|---|---|
| BFG | 16.50 |   |   |
| C | 26.50 | 30.67 |   |
| E | 32.00 | 33.00 | 41.00 |

|   | ADBFG | C |
|---|---|---|
| C | 29.00 |   |
| E | 32.60 | 41.00 |

# Hierarchical clustering example (UPGMA)



|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

|   | A | BF | C | D | E |
|---|---|----|---|---|---|
| BF | 18.50 |   |   |   |   |
| C | 27.00 | 31.50 |   |   |   |
| D | 8.00 | 17.50 | 26.00 |   |   |
| E | 33.00 | 35.50 | 41.00 | 31.00 |   |
| G | 13.00 | 12.50 | 29.00 | 14.00 | 28.00 |

|   | AD | BF | C | E |
|---|----|----|---|---|
| BF | 18.00 |   |   |   |
| C | 26.50 | 31.50 |   |   |
| E | 32.00 | 35.50 | 41.00 |   |
| G | 13.50 | 12.50 | 29.00 | 28.00 |

|   | AD | BFG | C |
|---|----|-----|---|
| BFG | 16.50 |   |   |
| C | 26.50 | 30.67 |   |
| E | 32.00 | 33.00 | 41.00 |

|   | ADBFG | C |
|---|-------|---|
| C | 29.00 |   |
| E | 32.60 | 41.00 |

|   | ADBFGC |
|---|--------|
| E | 34.00 |

Figure: Dr. Richard Edwards

# Hierarchical clustering example (UPGMA)

|   | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| A |   |   |   |   |   |   |   |
| B | 19.00 |   |   |   |   |   |   |
| C | 27.00 | 31.00 |   |   |   |   |   |
| D | 8.00 | 18.00 | 26.00 |   |   |   |   |
| E | 33.00 | 36.00 | 41.00 | 31.00 |   |   |   |
| F | 18.00 | 1.00 | 32.00 | 17.00 | 35.00 |   |   |
| G | 13.00 | 13.00 | 29.00 | 14.00 | 28.00 | 12.00 |   |

|   | A | BF | C | D | E |
|---|---|---|---|---|---|
| BF | 18.50 |   |   |   |   |
| C | 27.00 | 31.50 |   |   |   |
| D | 8.00 | 17.50 | 26.00 |   |   |
| E | 33.00 | 35.50 | 41.00 | 31.00 |   |
| G | 13.00 | 12.50 | 29.00 | 14.00 | 28.00 |

|   | AD | BF | C | E |
|---|---|---|---|---|
| BF | 18.00 |   |   |   |
| C | 26.50 | 31.50 |   |   |
| E | 32.00 | 35.50 | 41.00 |   |
| G | 13.50 | 12.50 | 29.00 | 28.00 |

|   | AD | BFG | C |
|---|---|---|---|
| BFG | 16.50 |   |   |
| C | 26.50 | 30.67 |   |
| E | 32.00 | 33.00 | 41.00 |

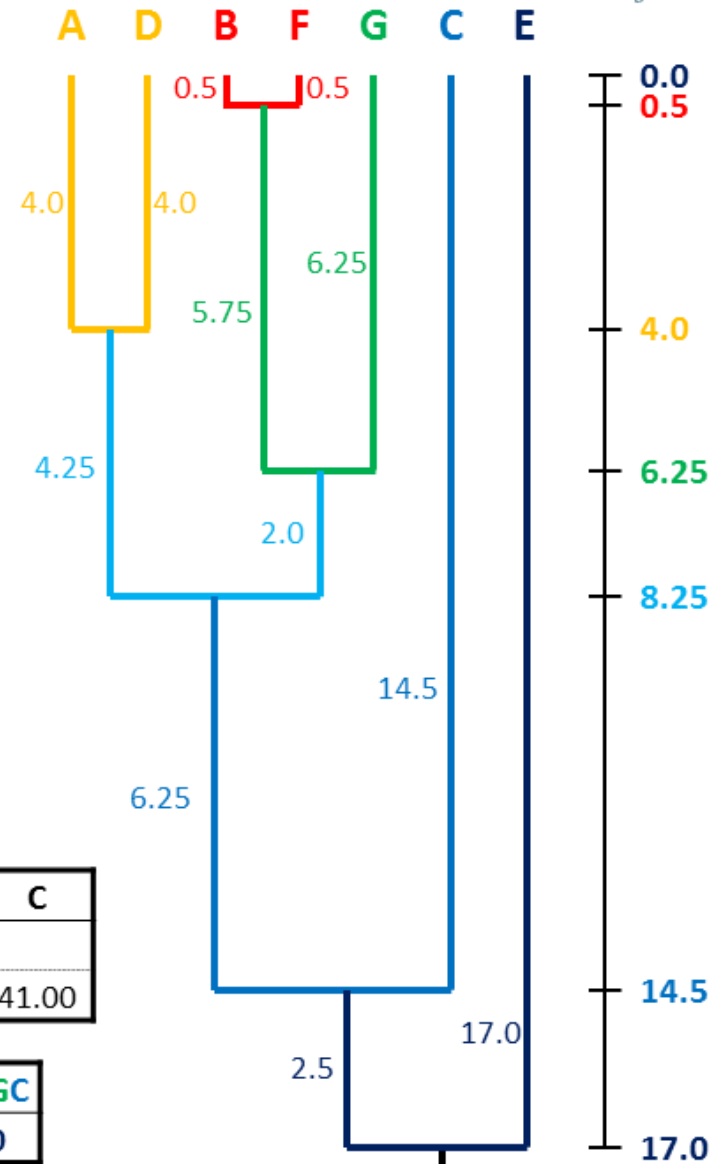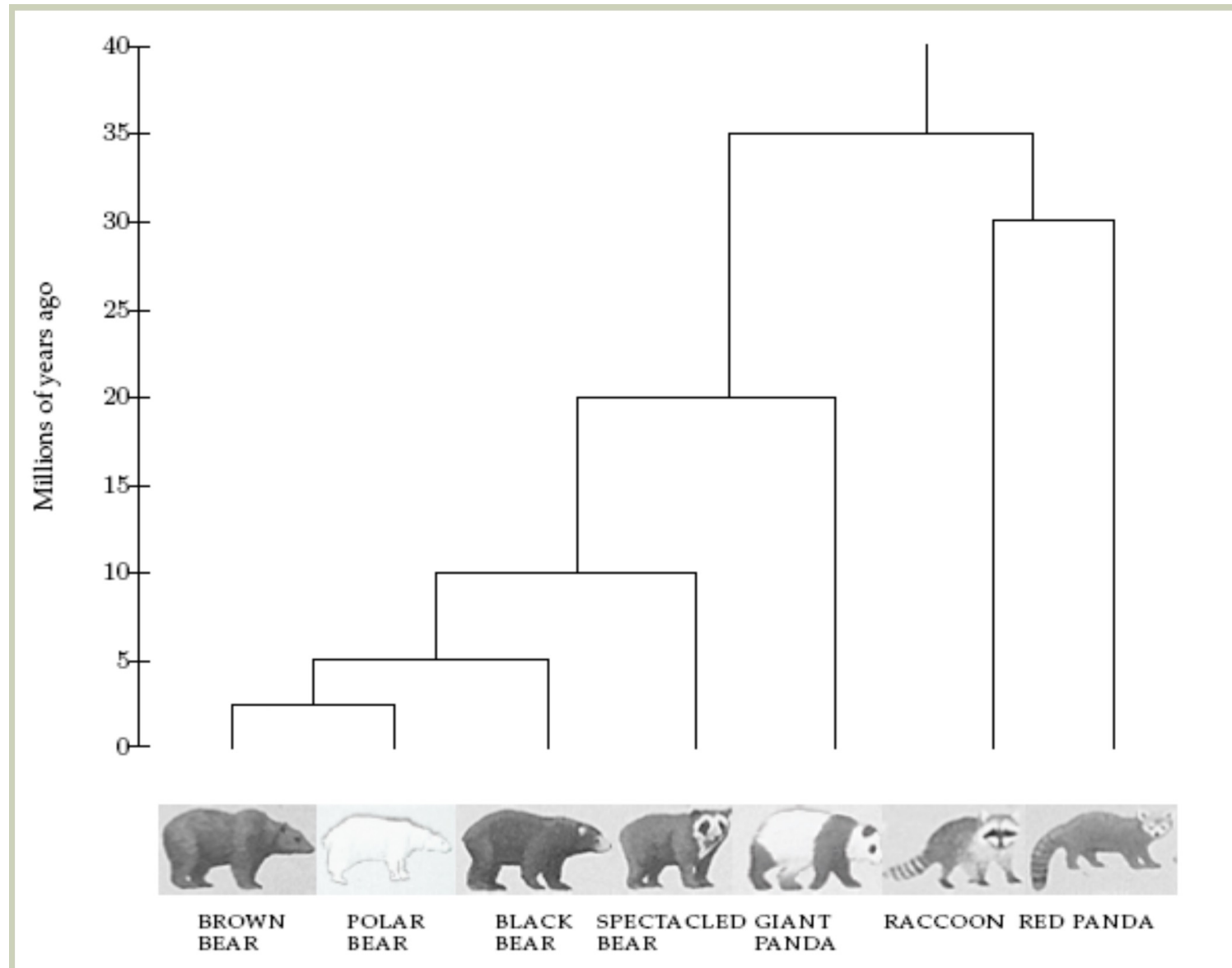|   | ADBFG | C |
|---|---|---|
| C | 29.00 |   |
| E | 32.60 | 41.00 |

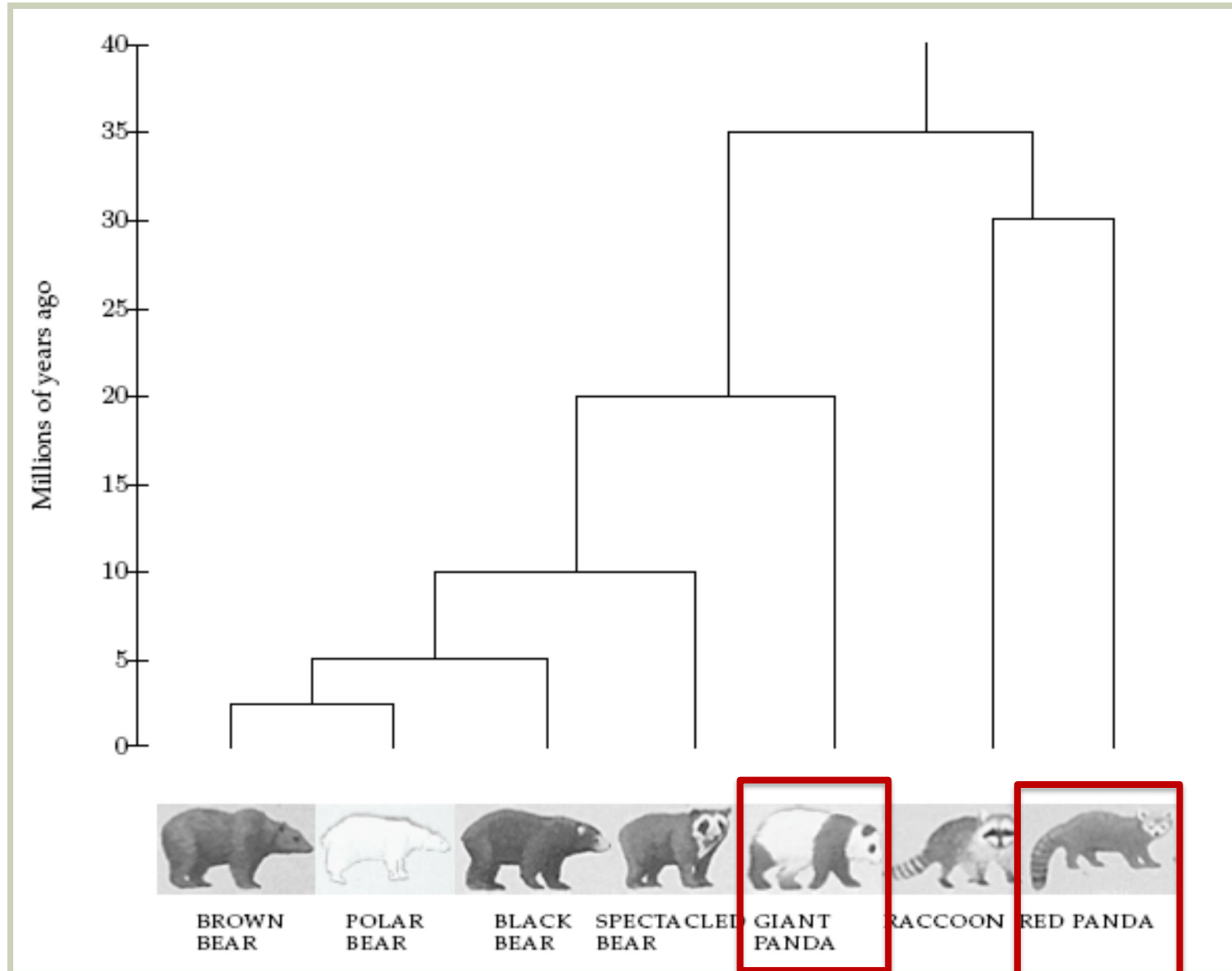|   | ADBFGC |
|---|---|
| E | 34.00 |



Figure: Dr. Richard Edwards

# Back to the pandas….

# Back to the pandas….



Credit: Ameet Soni

# Back to the pandas….



Credit: Ameet Soni