# Clustering

Kenneth (Kenny) Joseph

# Check your understanding

Consider trying k-means with different values of k. Which of the following graphs shows how the globally optimal heterogeneity changes for each value of k?

@_kenny_joseph

# How should we pick K then?

- The "Elbow rule"

@_kenny_joseph
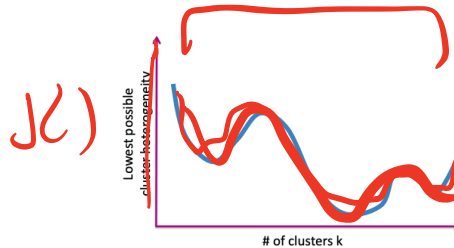
# Intuition: Local Minima, simple example
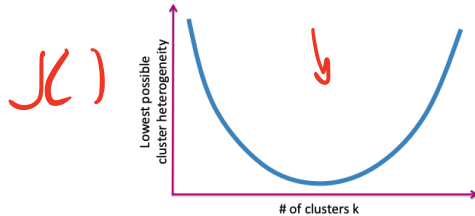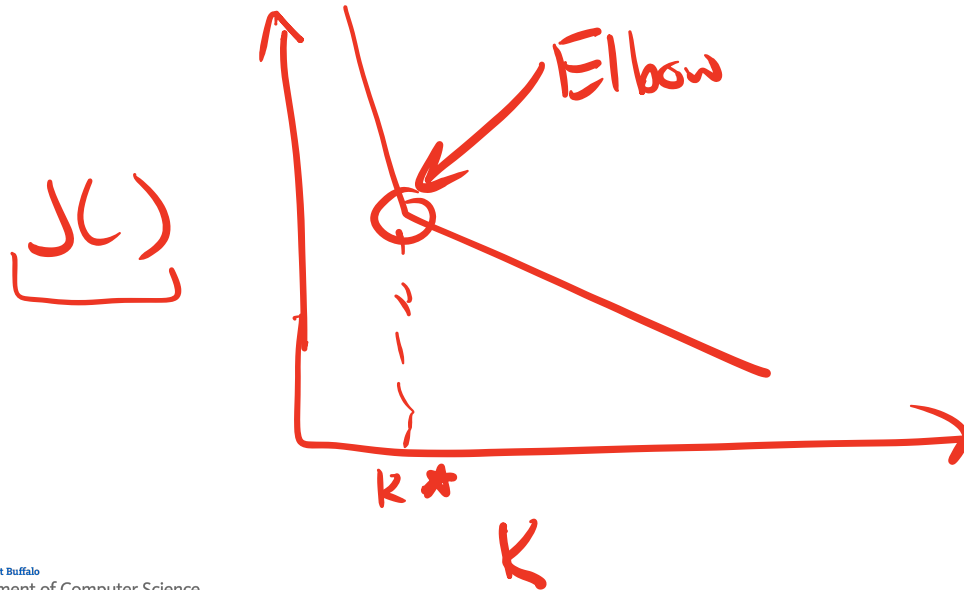
University at Buffalo
Department of Computer Science
and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# How do we evaluate?

**What makes one of these better than the other?**

@_kenny_joseph

# Two sets of evaluation metrics

- When the clusters are **known**
  - Can use the standard approaches, e.g. precision/recall (how?)
    - **PA4!**

- Can use a variety of metrics
  - Mutual information-based scores
  - Entropy-based scores
  - ...
- But we usually cluster when we *don't know* the labels!!

@_kenny_joseph

# Two sets of evaluation metrics

- When the clusters are **known**
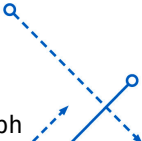  - Can use the standard approaches, e.g. precision/recall (how?)
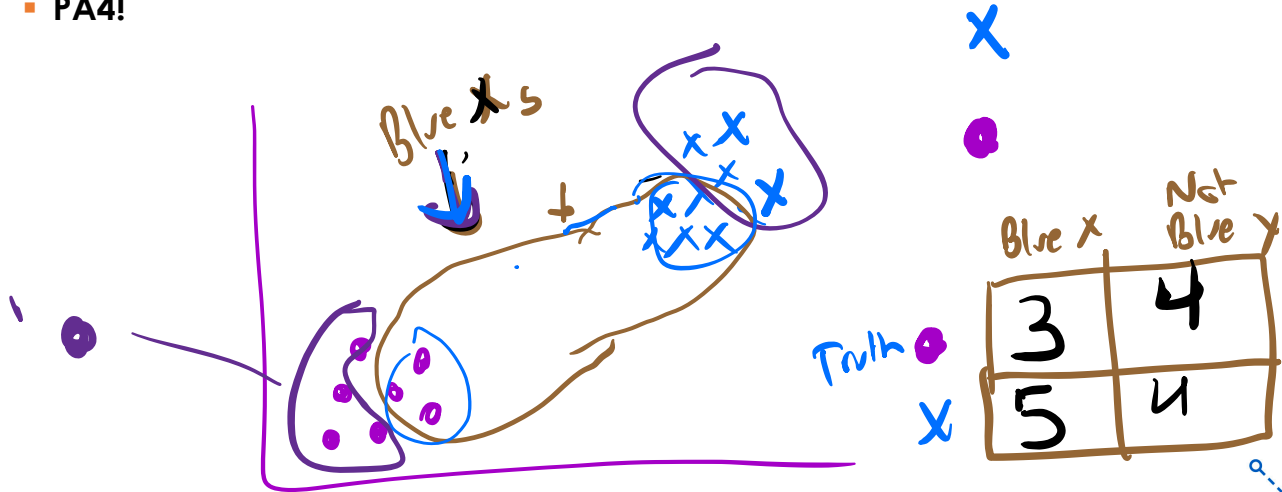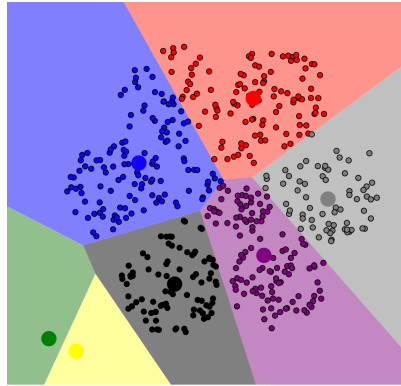    - **PA4!**

# How do we evaluate?

## What makes one of these better than the other?

University at Buffalo
Department of Computer Science
and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# Evaluation (con't)

**5 minutes: Come up with an evaluation *metric* that you could use to *quantify* your intuition. Give me a number, and how you computed it!**

@_kenny_joseph

# One Evaluation Metric – Silhouette Score

- **a**: The mean distance between a sample and all other points in the same class.
- **b**: The mean distance between a sample and all other points in the *next nearest cluster*.

The Silhouette Coefficient *s* for a single sample is then given as:

$$s = \frac{b - a}{max(a, b)}$$

*is defined*

*for a single point!* → *for dataset, just take average.*

University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences

@_kenny_joseph

# One Evaluation Metric – Silhouette Score

- **a**: The mean distance between a sample and all other points in the same class.
- **b**: The mean distance between a sample and all other points in the *next nearest cluster*.

The Silhouette Coefficient *s* for a single sample is then given as:

$$s = \frac{b - a}{max(a, b)}$$

sil score for ● ?

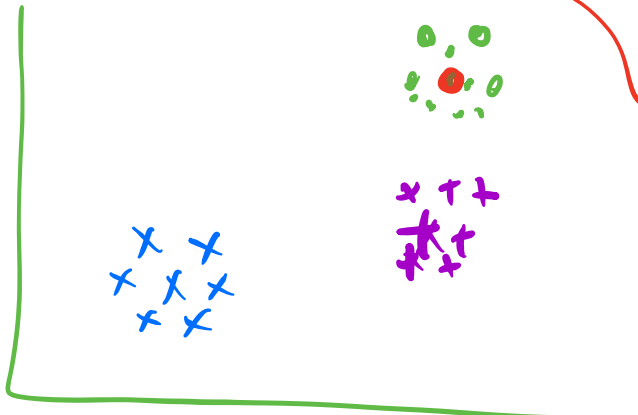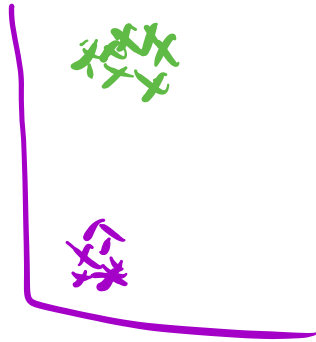a : 1. compute distance from ● to all the ● in its cluster
2. take mean.

b 1. Find nearest cluster: +
2. Distance from ● to all +
3. mean

11

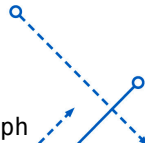@_kenny_joseph

# Code Demo

Intracluster distance should be small

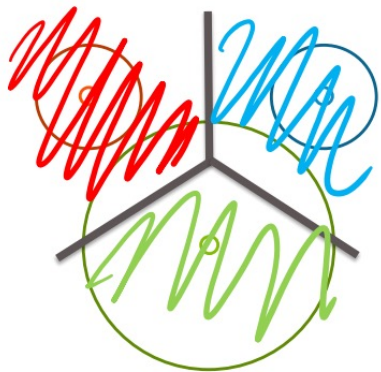Intercluster distance should be large

@_kenny_joseph

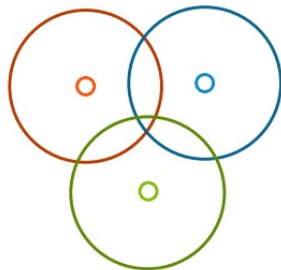# Kmeans Drawbacks: Difficulties w/ high dimensional data

- Full details: Section 3.5, CIML
- Intuition:
  - In high dimensions, distances start to become "more equal" (the variance of the distribution of distances across all points converges to a single number)
  - That's bad, because all kmeans does is work with distances between centers and points!
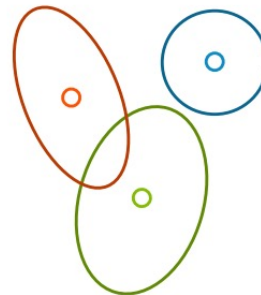  - Luckily, it's not all that bad, because points are not distributed uniformly,

@_kenny_joseph

# More Drawbacks to K-means



disparate cluster sizes

overlapping clusters

different shaped/oriented clusters

@_kenny_joseph

# A different approach: (Gaussian) mixture modeling

- Details in notebook...