



#### 10-601 Introduction to Machine Learning

Machine Learning Department School of Computer Science Carnegie Mellon University

# Clustering (K-Means)

#### **Clustering Readings:**

Murphy 25.5 Bishop 12.1, 12.3 HTF 14.3.0 Mitchell -- Matt Gormley Lecture 15 March 8, 2017

#### Reminders

- Homework 5: Readings / Application of ML
  - Release: Wed, Mar. 08
  - Due: Wed, Mar. 22 at 11:59pm

#### Outline

- Clustering: Motivation / Applications
- Optimization Background
  - Coordinate Descent
  - Block Coordinate Descent
- Clustering
  - Inputs and Outputs
  - Objective-based Clustering
- K-Means
  - K-Means Objective
  - Computational Complexity
  - K-Means Algorithm / Lloyd's Method
- K-Means Initialization
  - Random
  - Farthest Point
  - K-Means++

#### Clustering, Informal Goals

**Goal:** Automatically partition unlabeled data into groups of similar datapoints.

Question: When and why would we want to do this?

#### **Useful for:**

- Automatically organizing data.
- Understanding hidden structure in data.
- Preprocessing for further analysis.
  - Representing high-dimensional data in a low-dimensional space (e.g., for visualization purposes).

#### Applications (Clustering comes up everywhere...)

Cluster news articles or web pages or search results by topic.



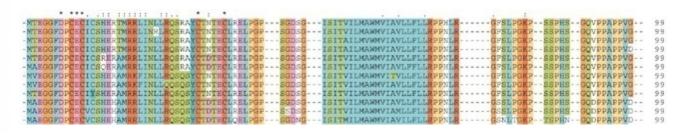




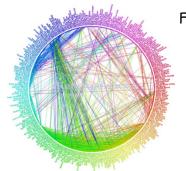


Cluster protein sequences by function or genes according to expression

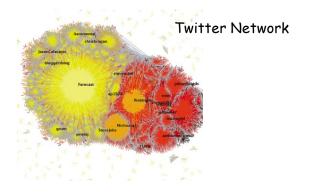
profile.



Cluster users of social networks by interest (community detection).



Facebook network



#### Applications (Clustering comes up everywhere...)

• Cluster customers according to purchase history.





Cluster galaxies or nearby stars (e.g. Sloan Digital Sky Survey)



And many many more applications....

# **Optimization Background**

- Coordinate Descent
- Block Coordinate Descent

# Clustering

- Inputs and Outputs
- Objective-based Clustering

#### **K-Means**

- K-Means Objective
- Computational Complexity
- K-Means Algorithm / Lloyd's Method

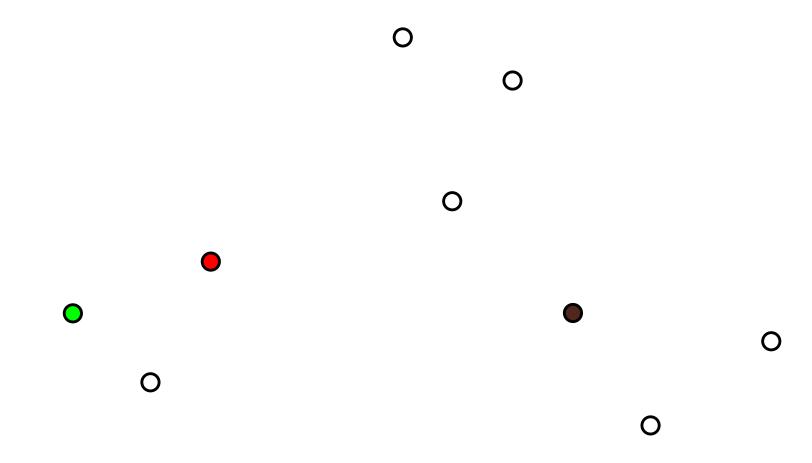
#### **K-Means Initialization**

- Random
- Furthest Traversal
- K-Means++

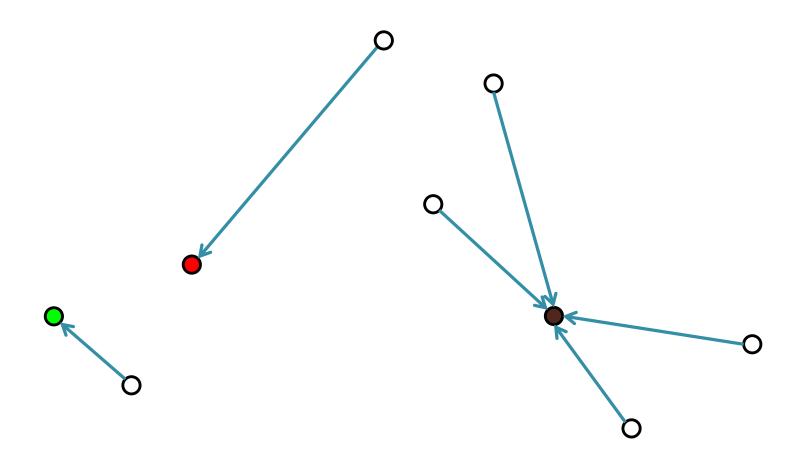
Example: Given a set of datapoints

			0					
					0			
				0				
		0						
0						0		C
	0							
							0	

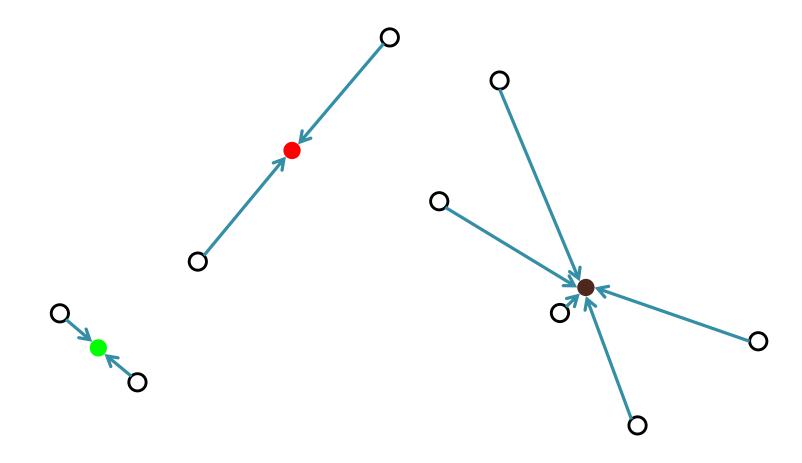
Select initial centers at random



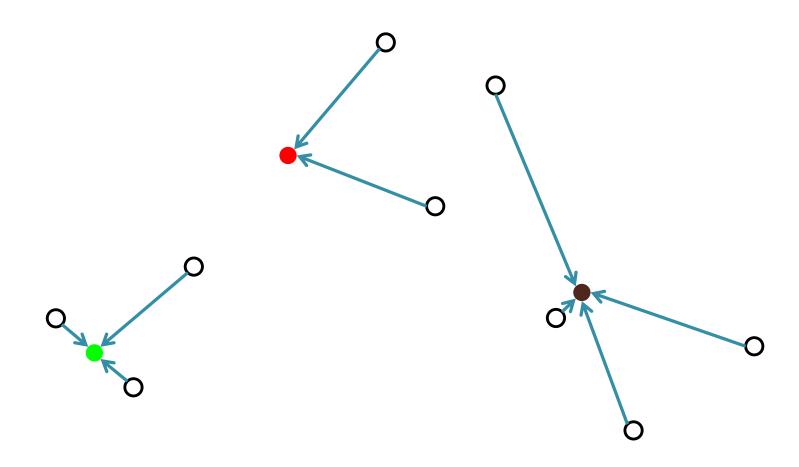
Assign each point to its nearest center



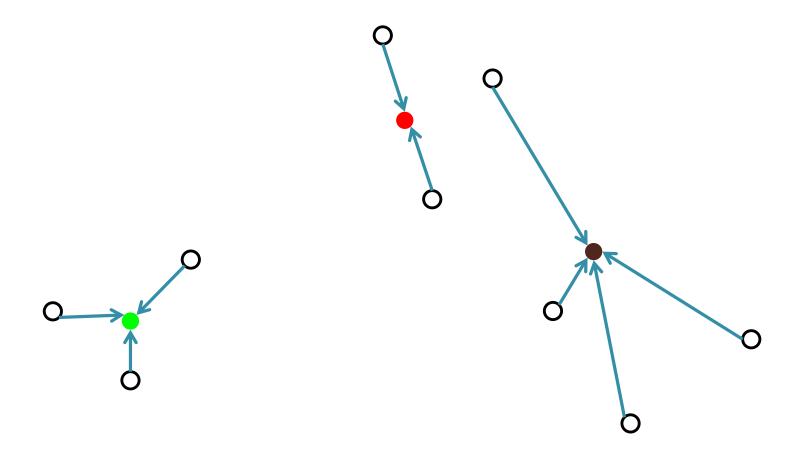
Recompute optimal centers given a fixed clustering



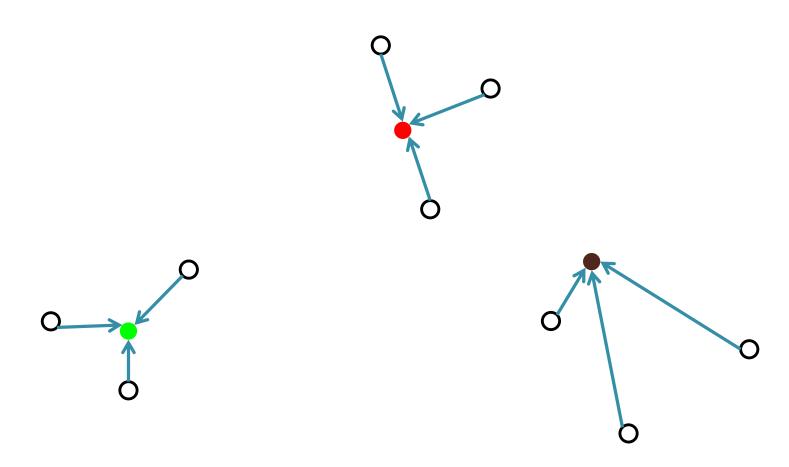
Assign each point to its nearest center



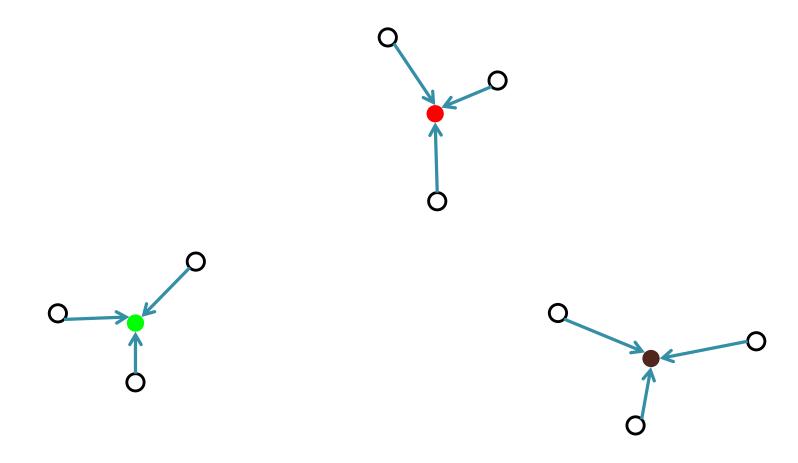
Recompute optimal centers given a fixed clustering



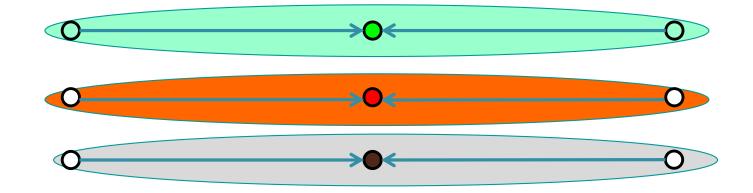
Assign each point to its nearest center



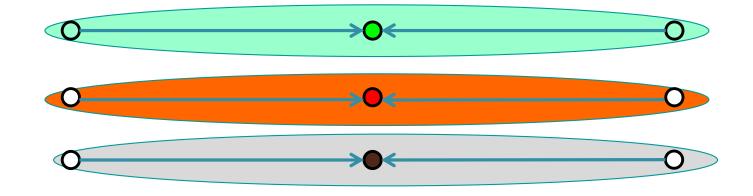
Recompute optimal centers given a fixed clustering



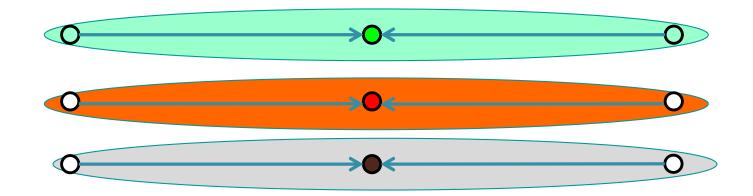
Get a good quality solution in this example.



It always converges, but it may converge at a local optimum that is different from the global optimum, and in fact could be arbitrarily worse in terms of its score.

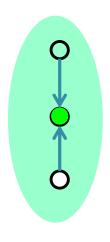


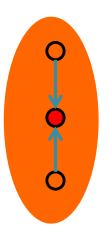
Local optimum: every point is assigned to its nearest center and every center is the mean value of its points.

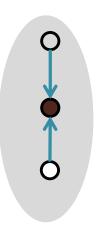


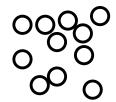
.It is arbitrarily worse than optimum solution....

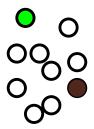


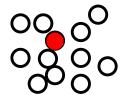




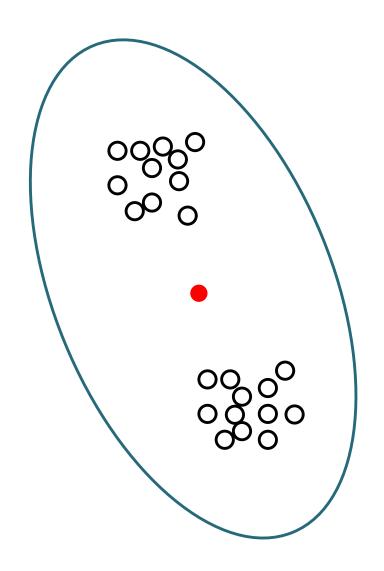


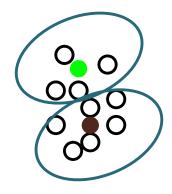






This bad performance, can happen even with well separated Gaussian clusters.





This bad performance, can happen even with well separated Gaussian clusters.

Some Gaussian are combined.....



If we do random initialization, as k increases, it becomes more likely
we won't have perfectly picked one center per Gaussian in our
initialization (so Lloyd's method will output a bad solution).

- For k equal-sized Gaussians, Pr[each initial center is in a different Gaussian]  $\approx \frac{k!}{k^k} \approx \frac{1}{e^k}$
- Becomes unlikely as k gets large.

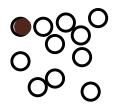
# Another Initialization Idea: Furthest Point Heuristic

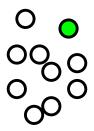
Choose c<sub>1</sub> arbitrarily (or at random).

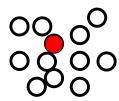
- For j = 2, ..., k
  - Pick  $c_j$  among datapoints  $x^1, x^2, ..., x^n$  that is farthest from previously chosen  $c_1, c_2, ..., c_{j-1}$

Fixes the Gaussian problem. But it can be thrown off by outliers....

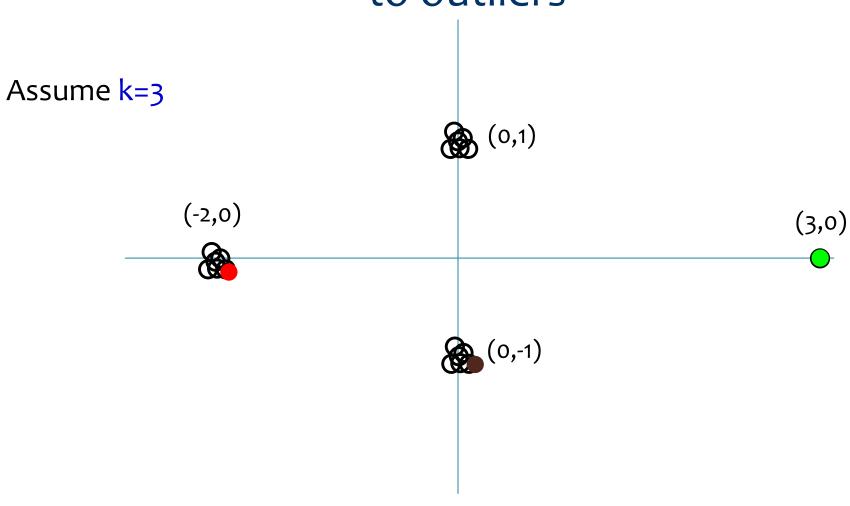
# Furthest point heuristic does well on previous example



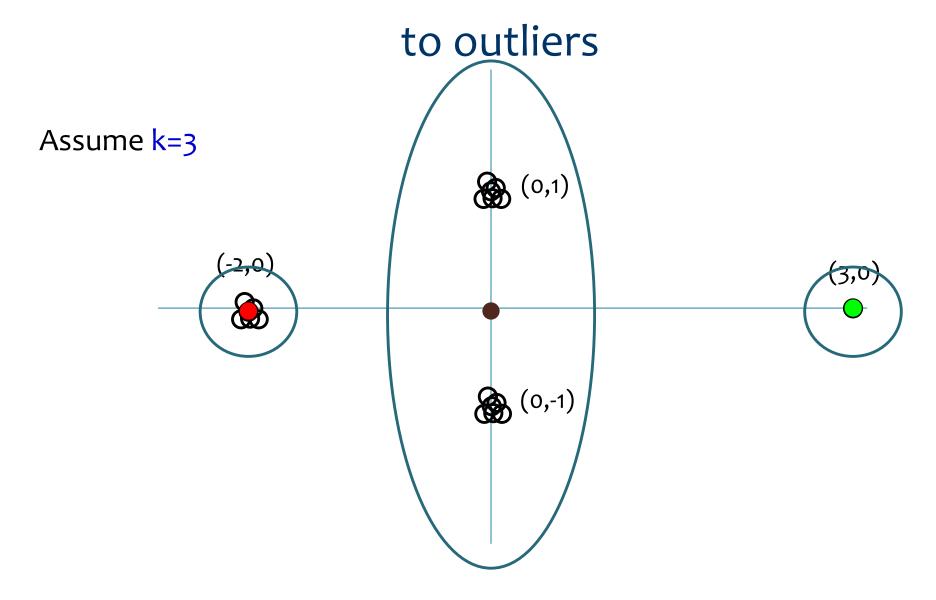




# Furthest point initialization heuristic sensitive to outliers



#### Furthest point initialization heuristic sensitive



#### K-means++ Initialization: D<sup>2</sup> sampling [AV07]

- Interpolate between random and furthest point initialization
- Let D(x) be the distance between a point x and its nearest center. Chose the next center proportional to  $D^2(x)$ .
  - Choose c<sub>1</sub> at random.
  - For j = 2, ..., k
    - Pick  $c_j$  among  $x^1, x^2, ..., x^n$  according to the distribution

$$Pr(c_j = x^i) \propto \min_{j' < j} \left| \left| x^i - c_{j'} \right| \right|^2 D^2(x^i)$$

**Theorem:** K-means++ always attains an O(log k) approximation to optimal k-means solution in expectation.

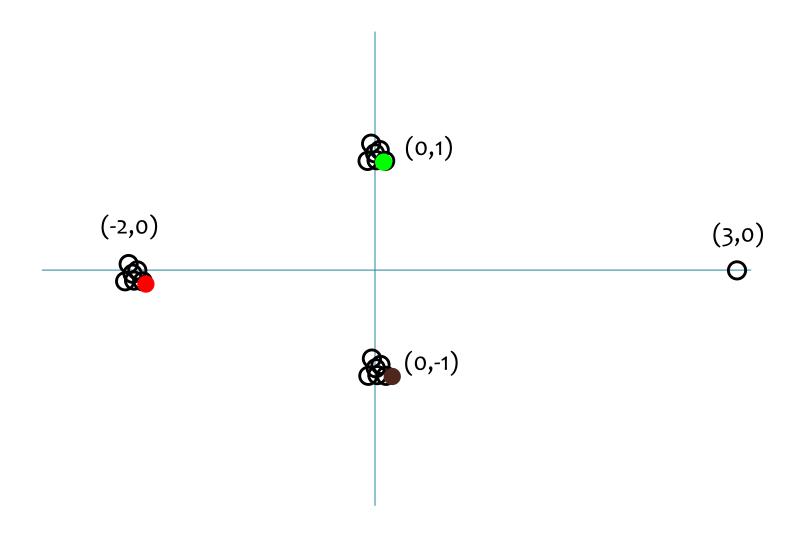
Running Lloyd's can only further improve the cost.

# K-means++ Idea: D<sup>2</sup> sampling

- Interpolate between random and furthest point initialization
- Let D(x) be the distance between a point x and its nearest center. Chose the next center proportional to  $D^{\alpha}(x)$ .
  - $\alpha = 0$ , random sampling
  - $\alpha = \infty$ , furthest point (Side note: it actually works well for k-center)
  - $\alpha = 2$ , k-means++

Side note:  $\alpha = 1$ , works well for k-median

#### K-means ++ Fix



#### K-means++/ Lloyd's Running Time

- K-means ++ initialization: O(nd) and one pass over data to select next center. So O(nkd) time in total.
- Lloyd's method

**Repeat** until there is no change in the cost.

- For each j:  $C_i \leftarrow \{x \in S \text{ whose closest center is } c_i\}$ 
  - For each j: c<sub>i</sub> ←mean of C<sub>i</sub>

Each round takes time O(nkd).

- Exponential # of rounds in the worst case [AVo7].
- Expected polynomial time in the smoothed analysis (non worst-case) model!

#### K-means++/ Lloyd's Summary

- Running Lloyd's can only further improve the cost.
- Exponential # of rounds in the worst case [AVo7].
- Expected polynomial time in the smoothed analysis model!
- Does well in practice.
- K-means++ always attains an O(log k) approximation to optimal k-means solution in expectation.

#### What value of k???

Heuristic: Find large gap between k -1-means cost and k-means cost.

• Hold-out validation/cross-validation on auxiliary task (e.g., supervised learning task).

Try hierarchical clustering.