

10-301/601: Introduction to Machine Learning Lecture 19: Clustering

Henry Chai

7/11/23

Front Matter

- Announcements
 - PA4 released 6/15, due 7/13 at 11:59 PM
- Recommended Readings
 - Murphy, [Chapters 25.5.1 - 25.5.2](#)
 - Daumé III, [Chapter 15: Unsupervised Learning](#)

Learning Paradigms

- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$
 - Regression - $y^{(n)} \in \mathbb{R}$
 - Classification - $y^{(n)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$
 - **Clustering**
 - Dimensionality reduction

Clustering

- Goal: split an unlabeled data set into groups or clusters of "similar" data points
- Use cases:
 - Organizing data
 - Discovering patterns or structure
 - Preprocessing for downstream machine learning tasks
- Applications:
 - customer grouping
 - determining address information (blocks?)
 - distinguishing categories of items from pictures/3D maps

Recall: Similarity for k NN

- Intuition: ~~predict the label of a data point to be the label of the “most similar” training point~~ two points are “similar” if the distance between them is small
- Euclidean distance: $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2$

Partition-Based Clustering

- Given a desired number of clusters, K , return a partition of the data set into K groups or clusters, $\{C_1, \dots, C_K\}$, that optimize some objective function
 1. What objective function should we optimize?
 2. How can we perform optimization in this setting?



Option A



Option B

Which do you prefer?

Which partition do you prefer?

Option
A

Option
B

General Recipe for Machine Learning

- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

Recipe for K-means

- Define a model and model parameters
 - Assume there are K clusters
 - Use the Euclidean distance
 - Parameters: cluster centers: μ_1, \dots, μ_K

- Write down an objective function

$$L(\mu_1, \dots, \mu_K, z^{(1)}, \dots, z^{(N)})$$
$$= \sum_{n=1}^N \|x^{(n)} - \mu_{z^{(n)}}\|_2^2$$

cluster assignments:
 $z^{(1)}, \dots, z^{(N)}$

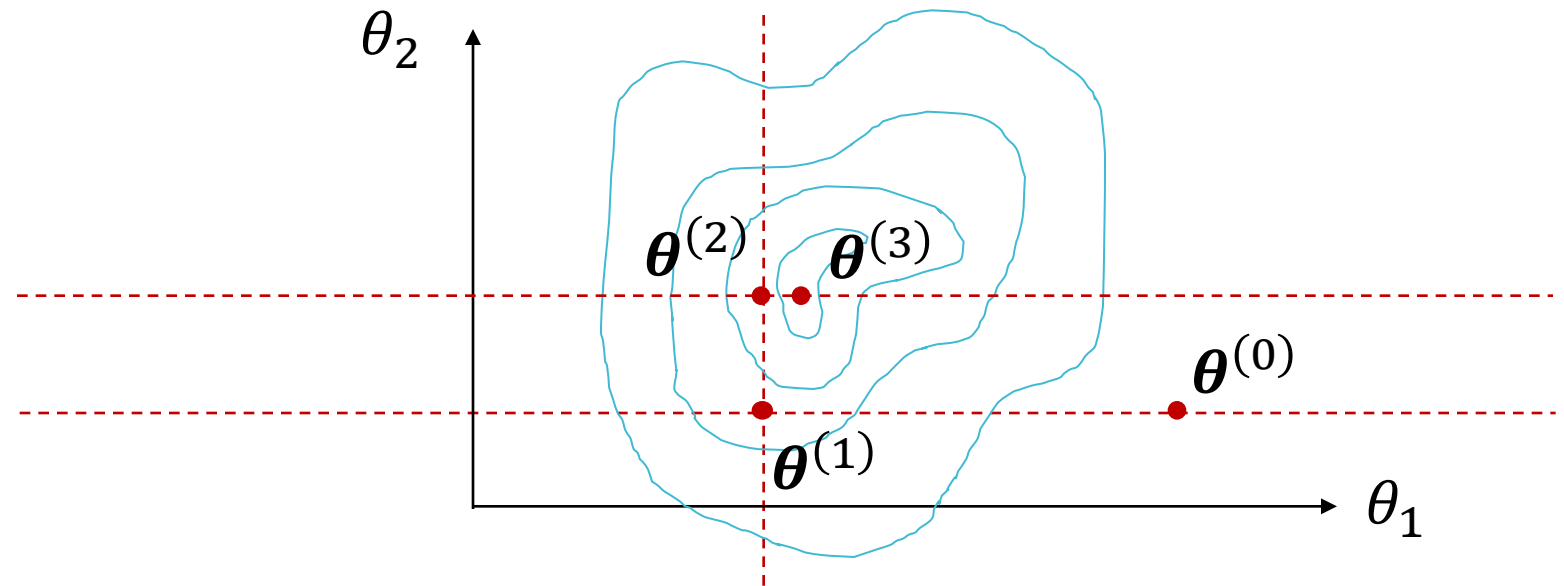
- Optimize the objective w.r.t. the model parameters
 - Use block coordinate descent

Coordinate Descent

- Goal: minimize some objective

$$\hat{\theta} = \operatorname{argmin} J(\theta)$$

- Idea: iteratively pick one variable and minimize the objective w.r.t. just that variable, *keeping all others fixed*.



Block Coordinate Descent

- Goal: minimize some objective

$$\hat{\alpha}, \hat{\beta} = \operatorname{argmin} J(\alpha, \beta)$$

- Idea: iteratively pick one *block* of variables (α or β) and minimize the objective w.r.t. that block, keeping the other(s) fixed.
 - Ideally, blocks should be the largest possible set of variables *that can be efficiently optimized simultaneously*

Optimizing the K -means objective

$$\hat{\mu}_1, \dots, \hat{\mu}_K, \hat{z}^{(1)}, \dots, \hat{z}^{(N)} = \operatorname{argmin} \sum_{n=1}^N \|x^{(n)} - \mu_{z^{(n)}}\|_2$$

- If μ_1, \dots, μ_K are fixed

$$\hat{z}^{(n)} = \operatorname{argmin}_{k \in \{1, \dots, K\}} \|x^{(n)} - \mu_k\|_2$$

- If $z^{(1)}, \dots, z^{(N)}$ are fixed

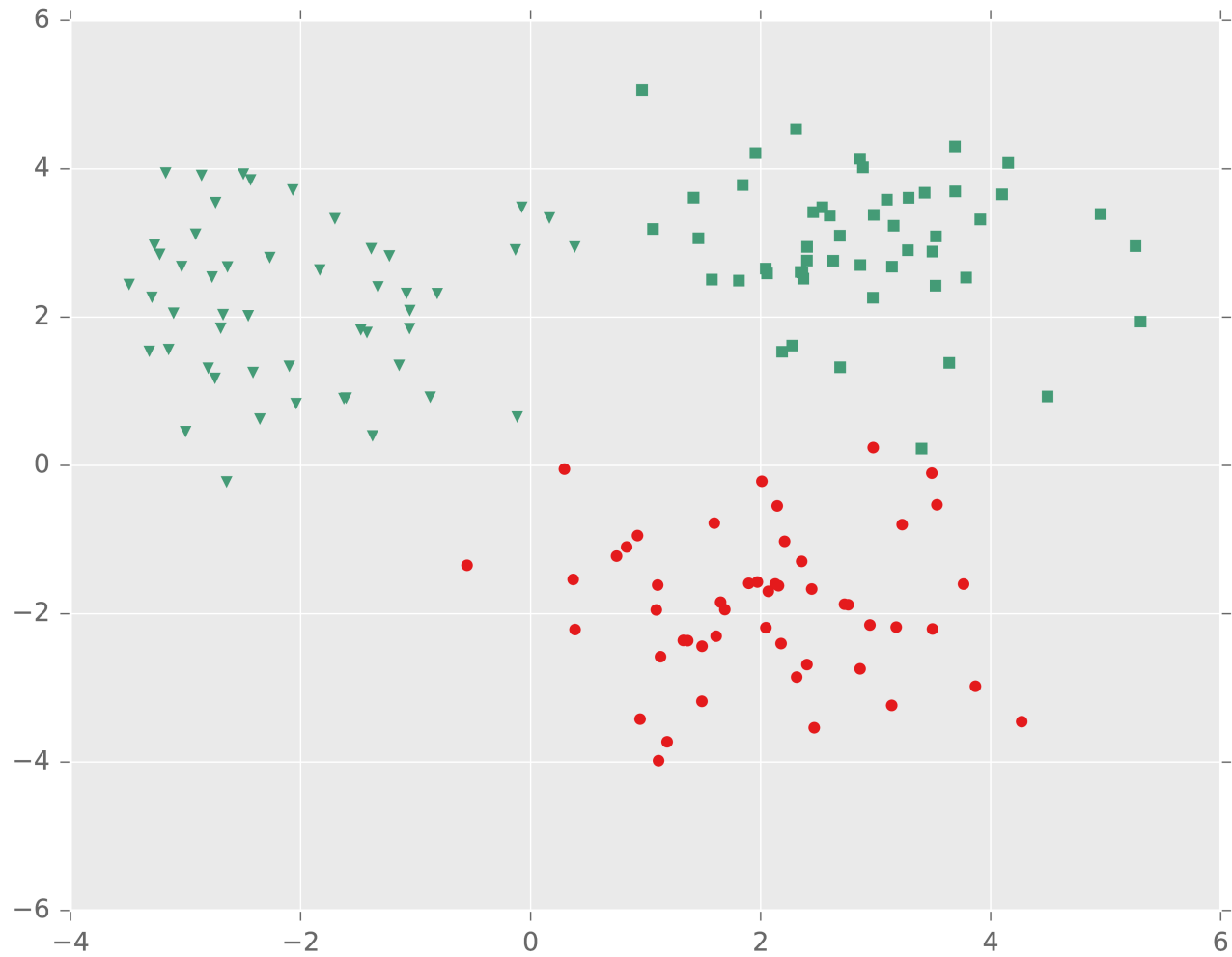
$$\hat{\mu}_k = \frac{1}{N_k} \sum_{n: z^{(n)} = k} x^{(n)}$$

where $N_k = \#$ of data points in cluster k

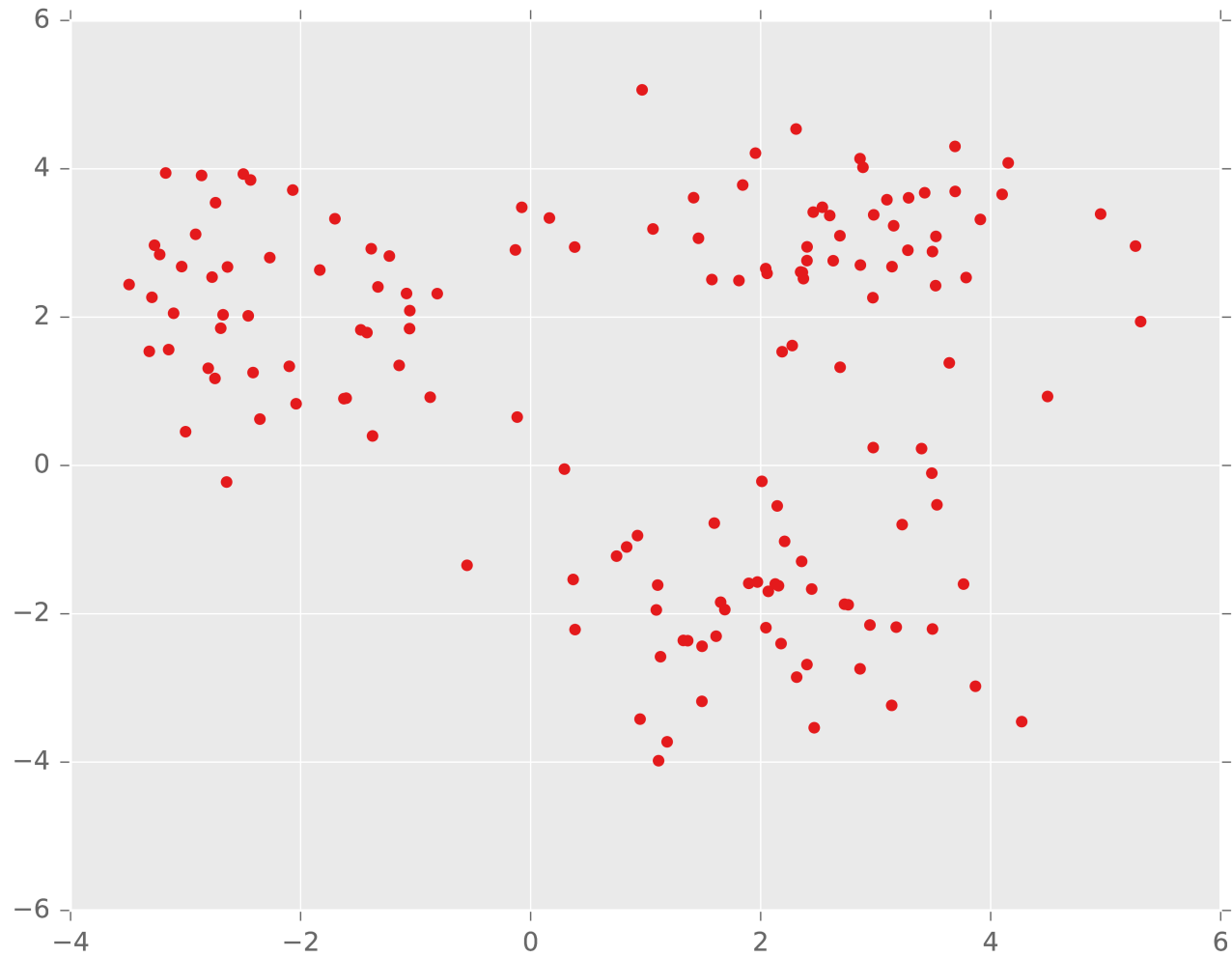
K-means Algorithm

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, K$
 1. Initialize cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$
 2. While NOT CONVERGED
 - a. Assign each data point to the cluster with the nearest cluster center:
$$z^{(n)} = \underset{k}{\operatorname{argmin}} \|\mathbf{x}^{(n)} - \boldsymbol{\mu}_k\|_2$$
 - b. Recompute the cluster centers:
$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n: z^{(n)}=k} \mathbf{x}^{(n)}$$
where N_k is the number of data points in cluster k
- Output: cluster centers $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ and cluster assignments $z^{(1)}, \dots, z^{(N)}$

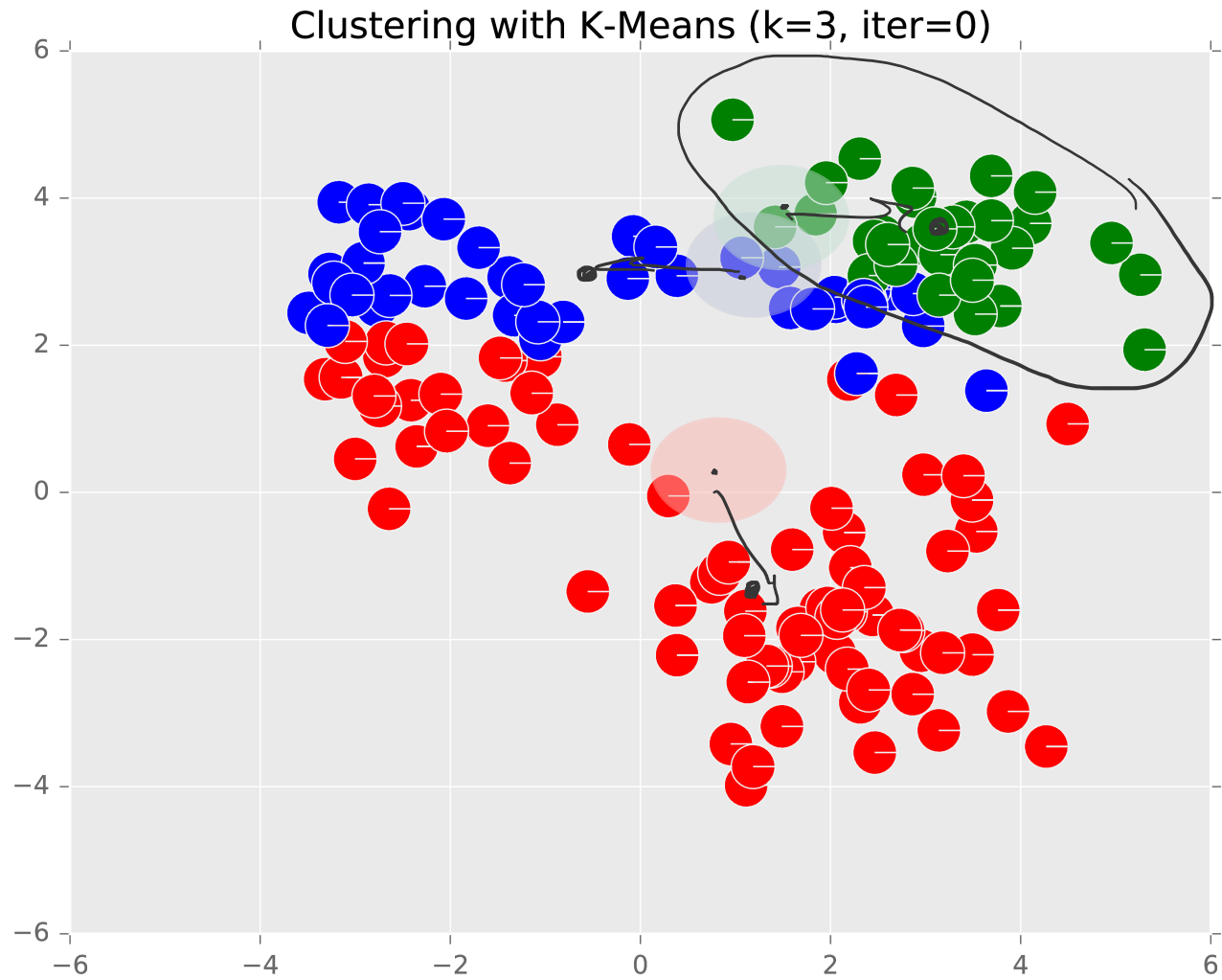
K -means: Example ($K = 3$)



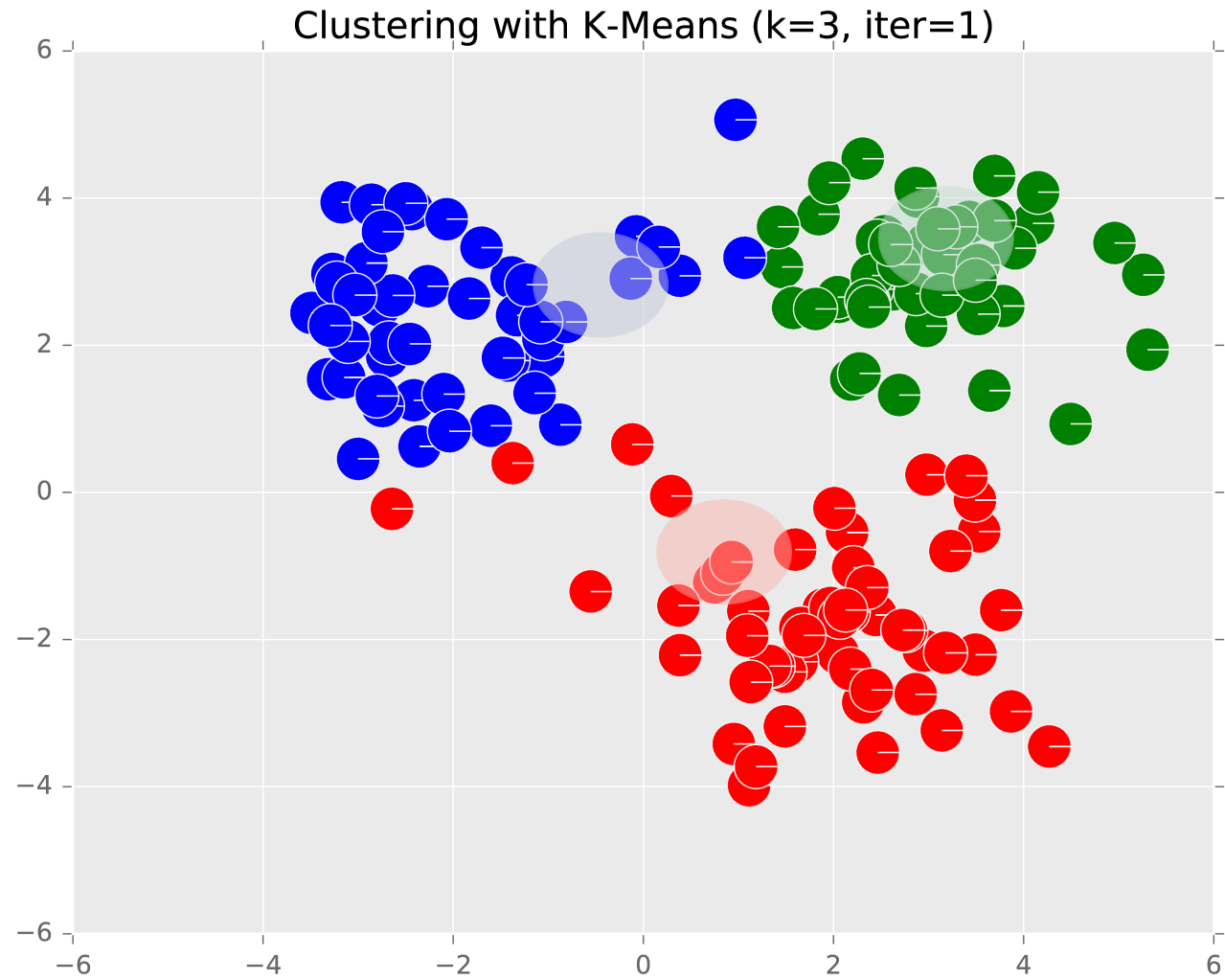
K -means: Example ($K = 3$)



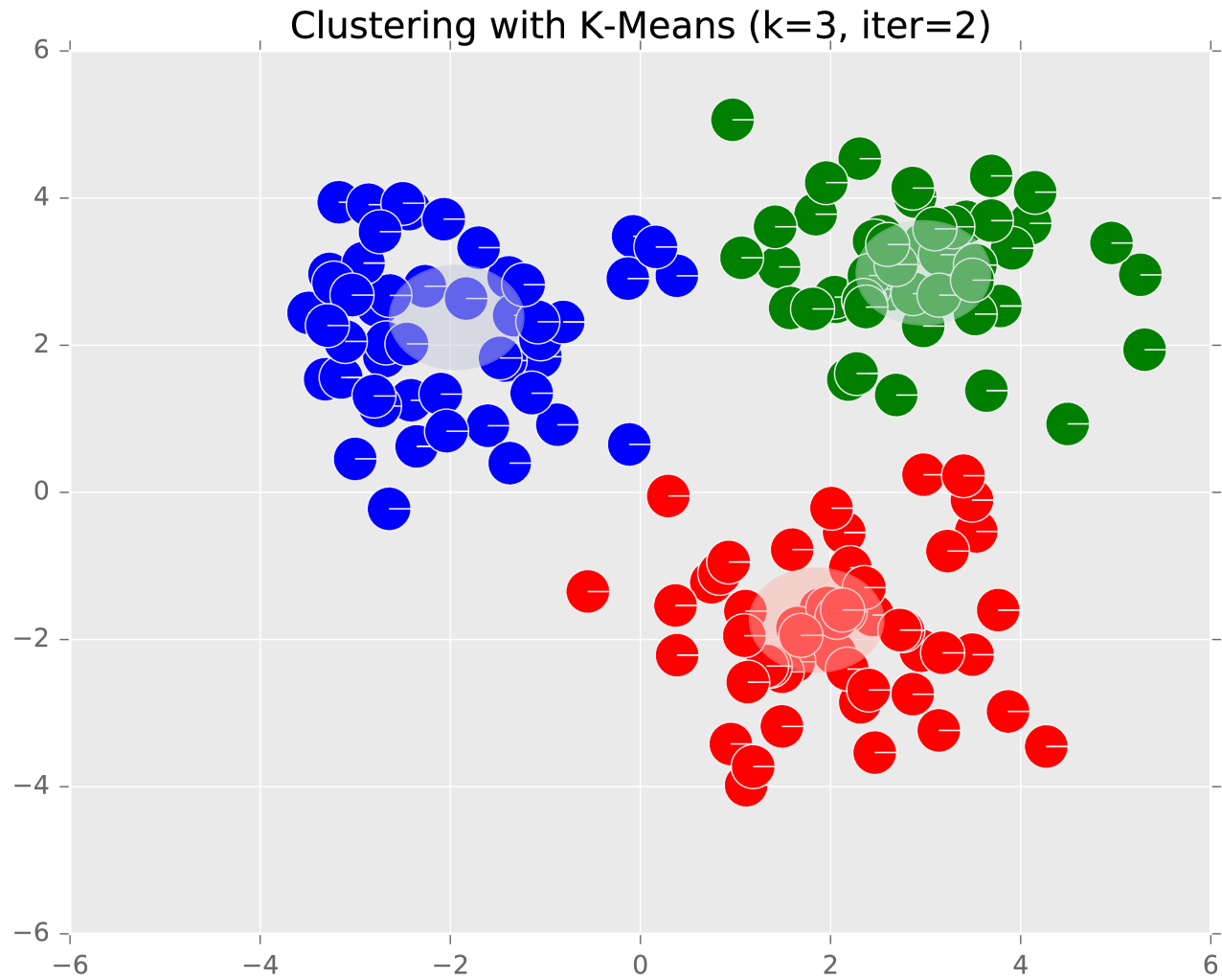
K -means: Example ($K = 3$)



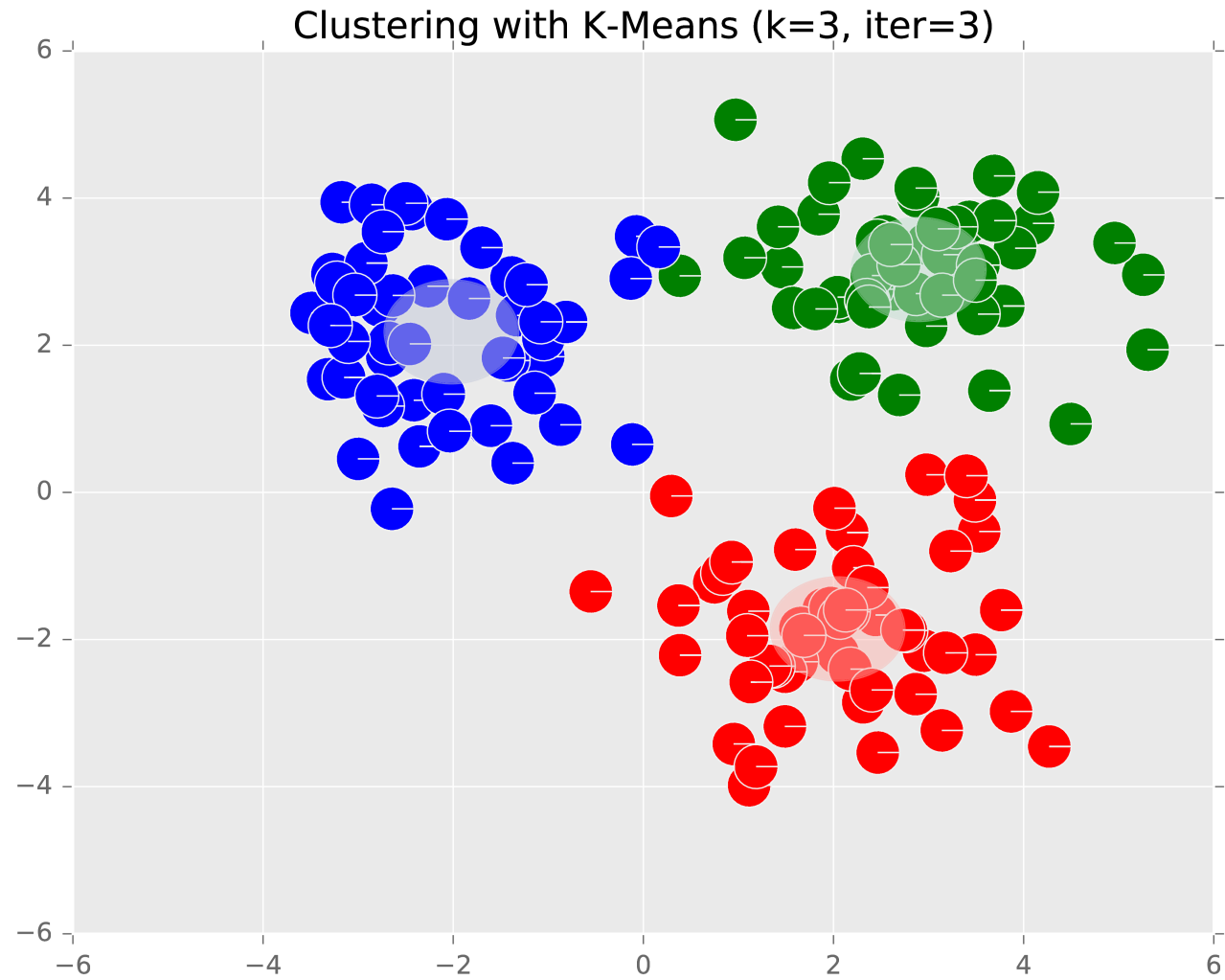
K-means:
Example
($K = 3$)



K-means:
Example
($K = 3$)



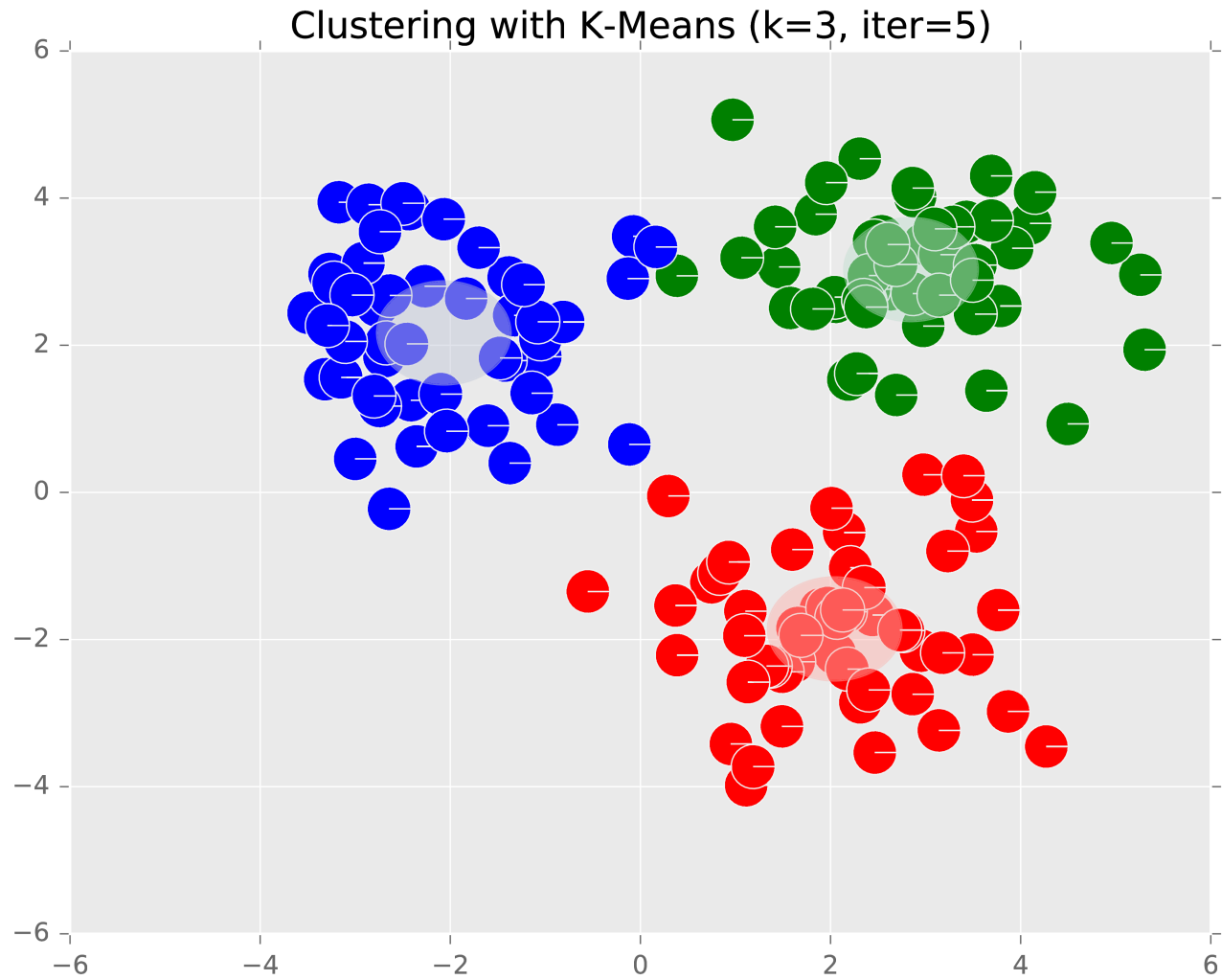
K-means:
Example
($K = 3$)



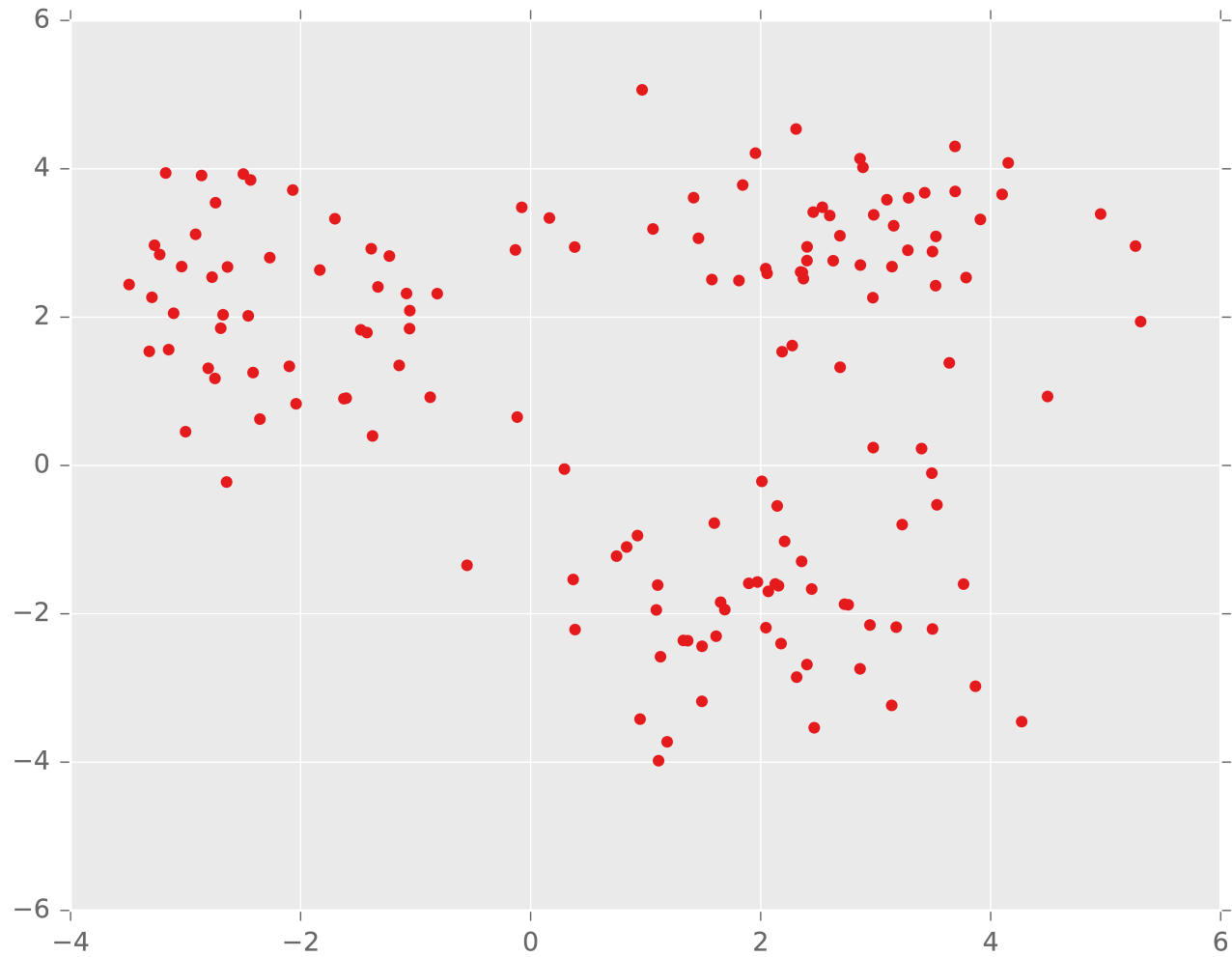
K-means:
Example
($K = 3$)



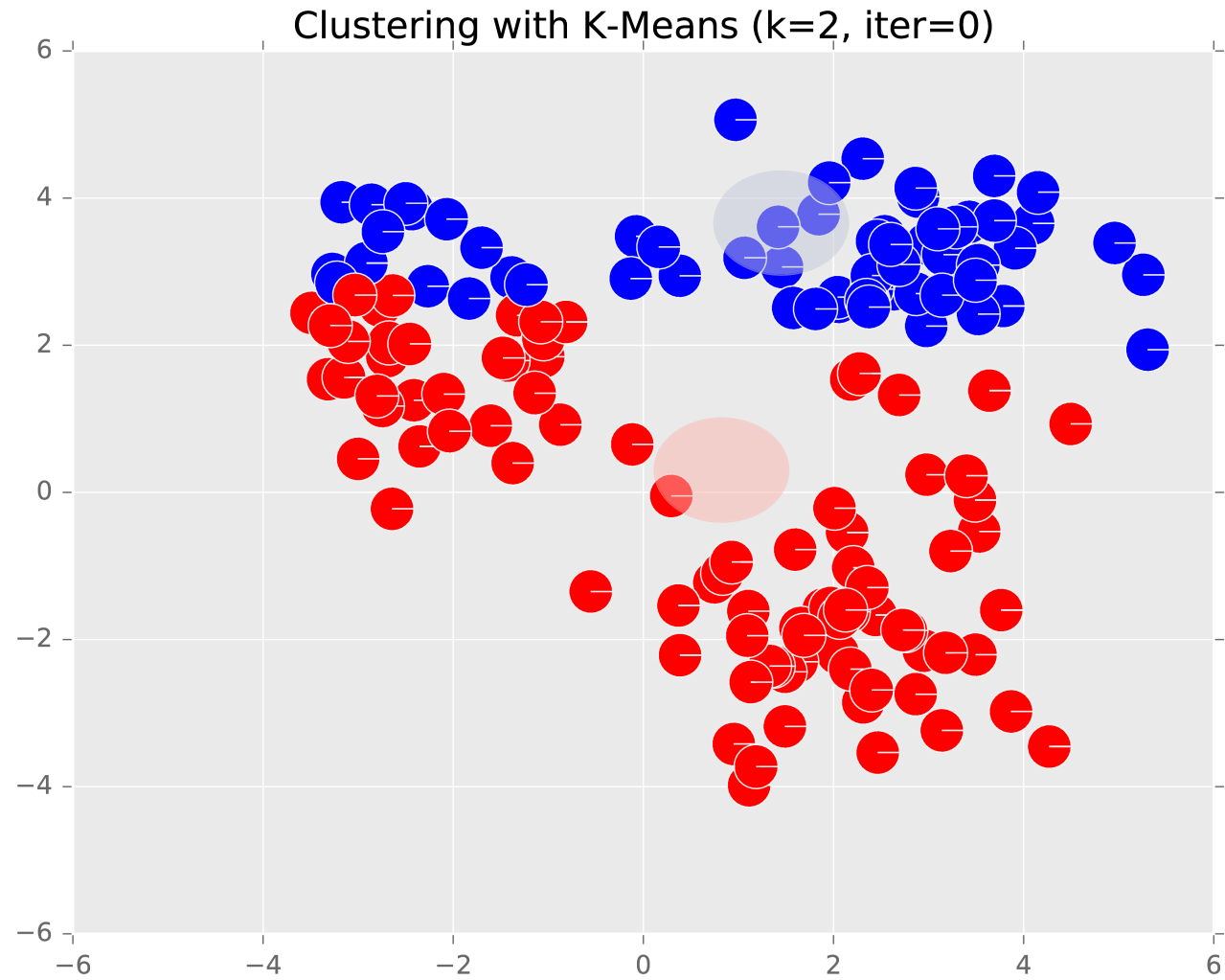
K-means:
Example
($K = 3$)



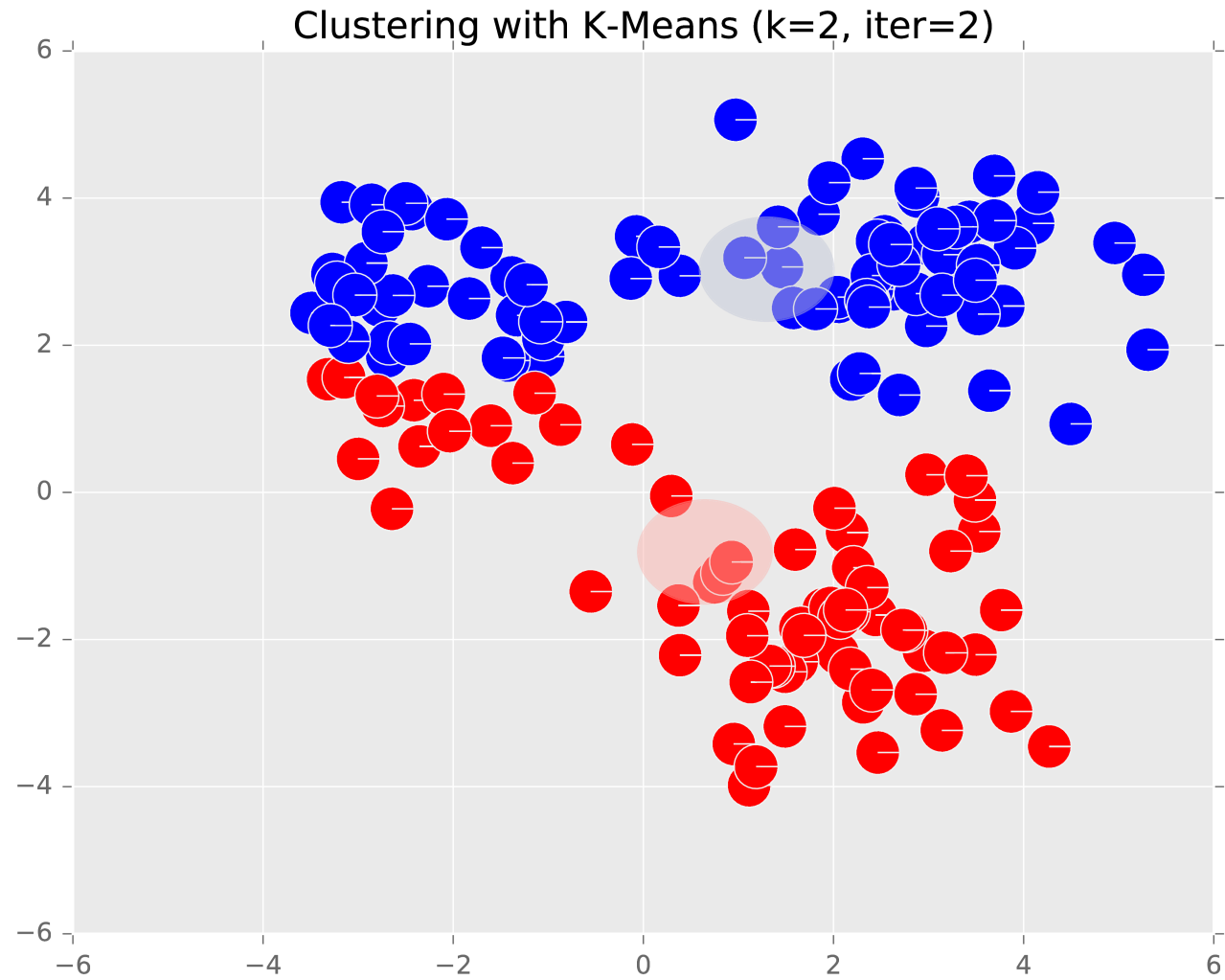
K -means: Example ($K = 2$)



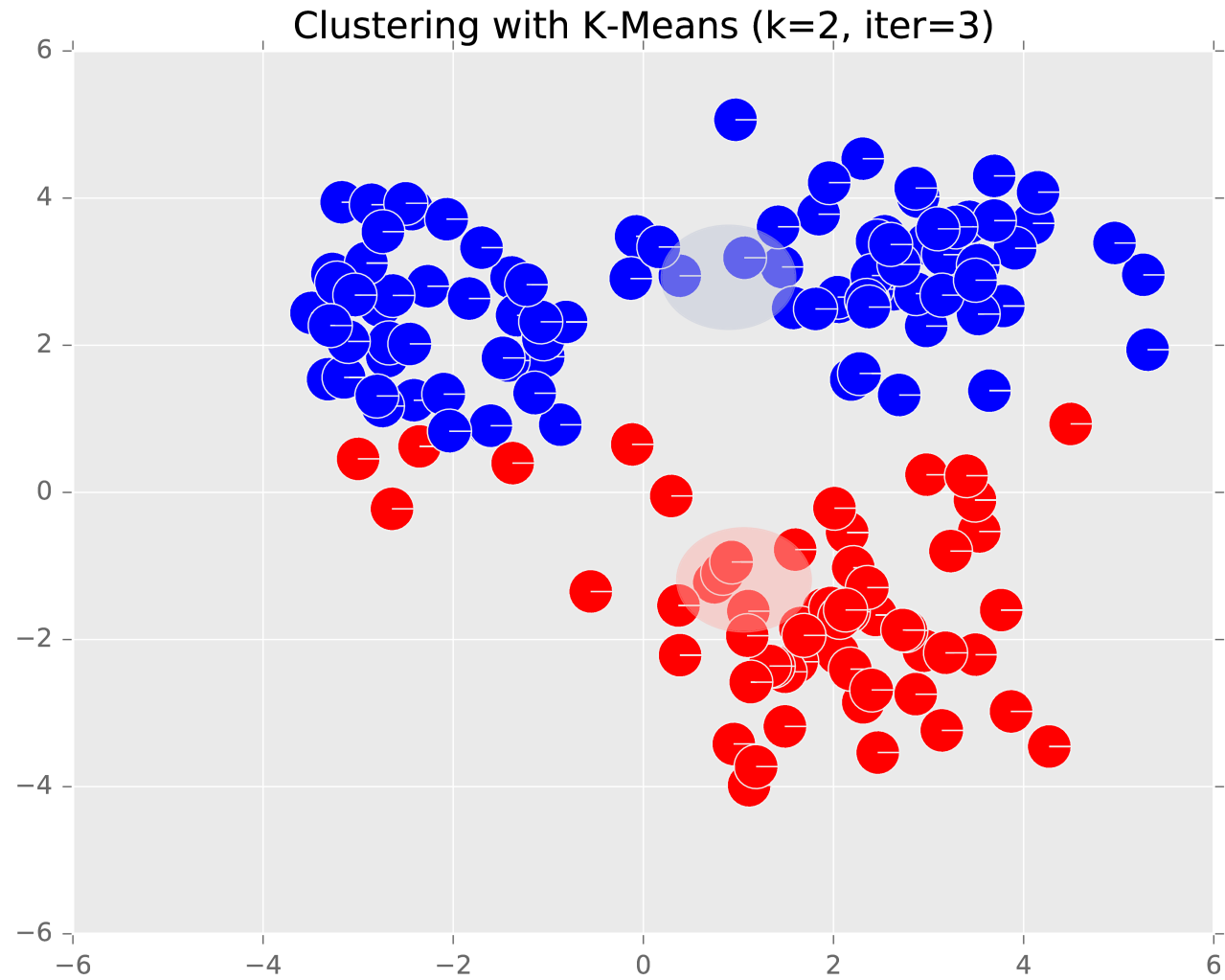
K-means:
Example
($K = 2$)



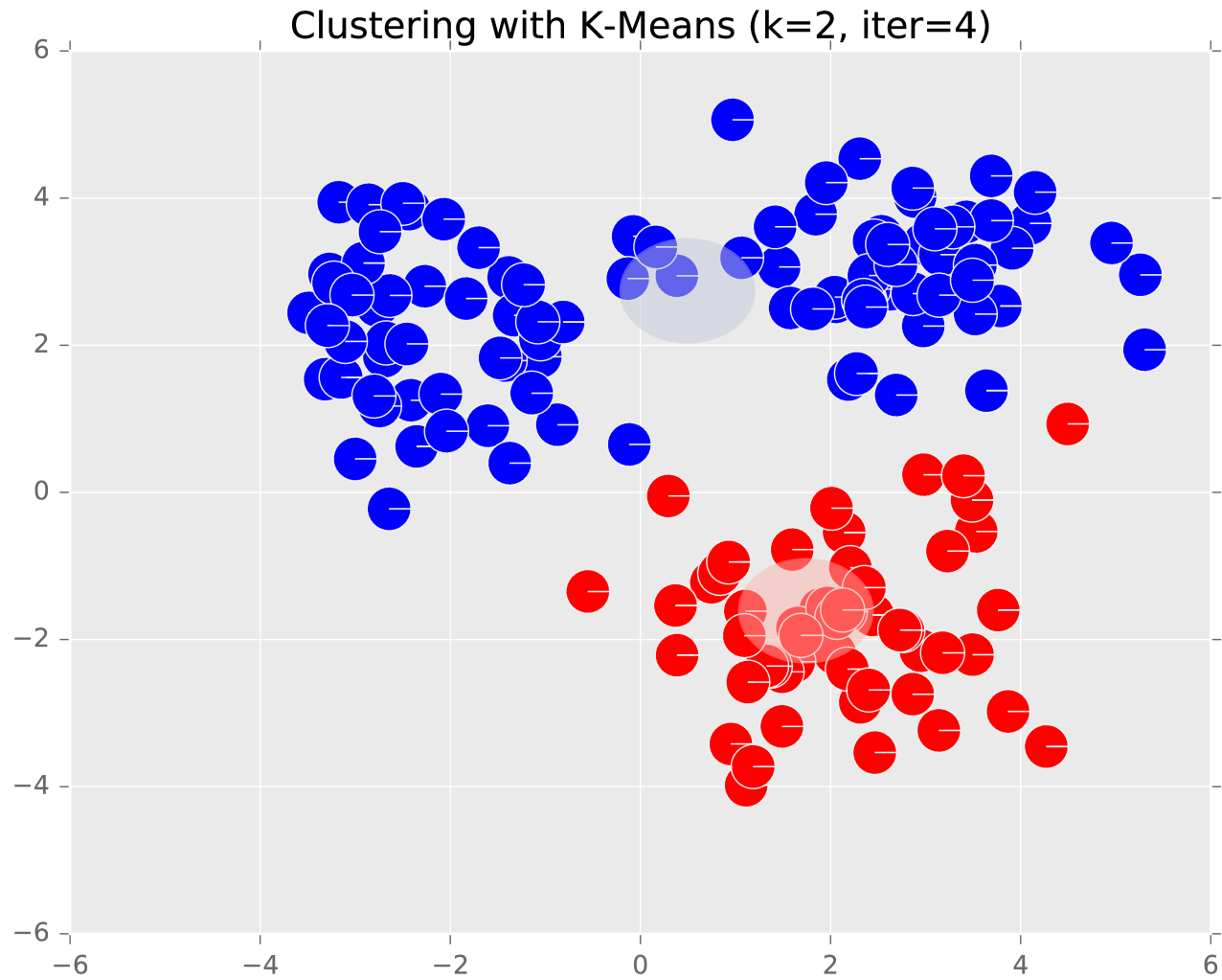
K-means:
Example
($K = 2$)



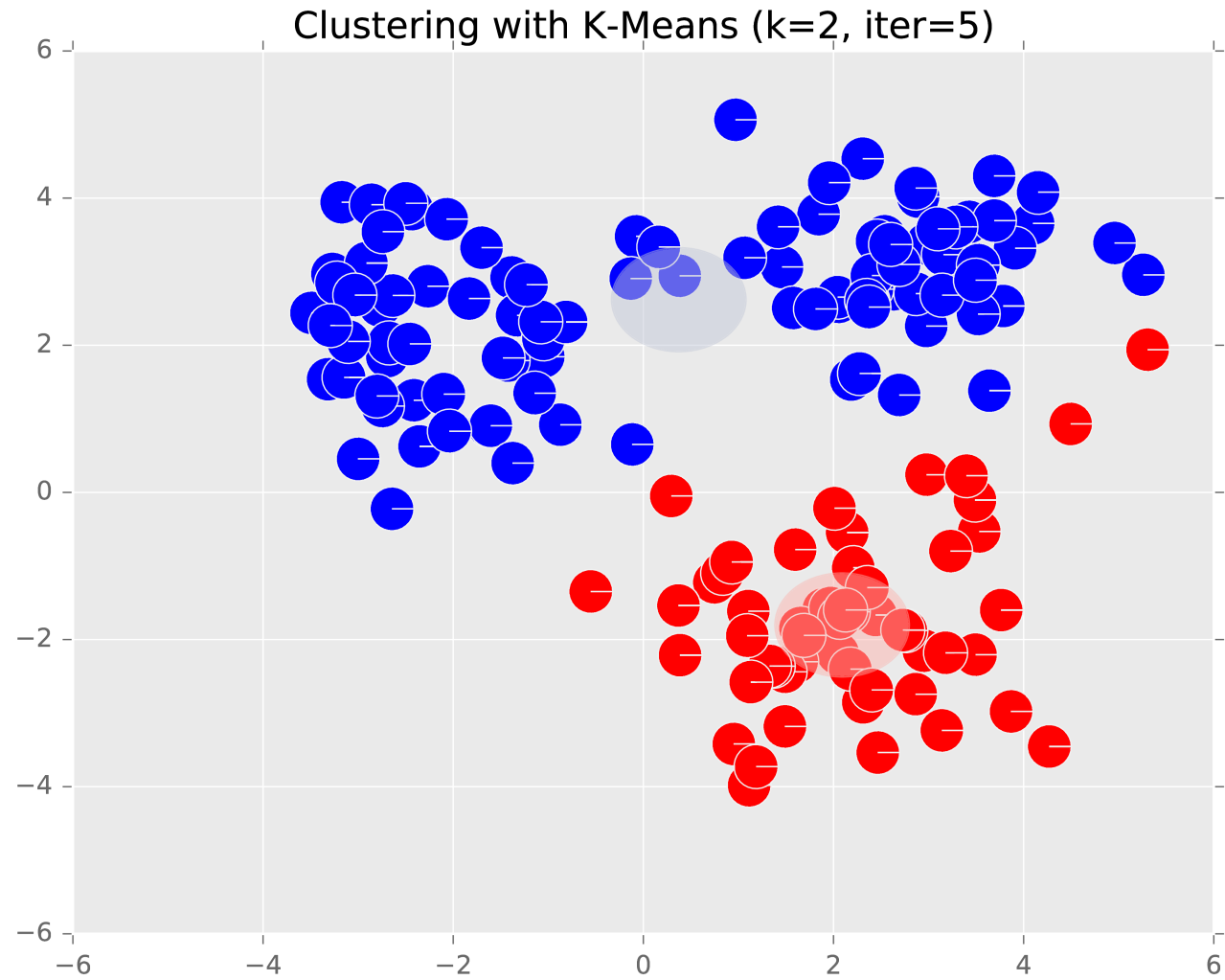
K-means:
Example
($K = 2$)



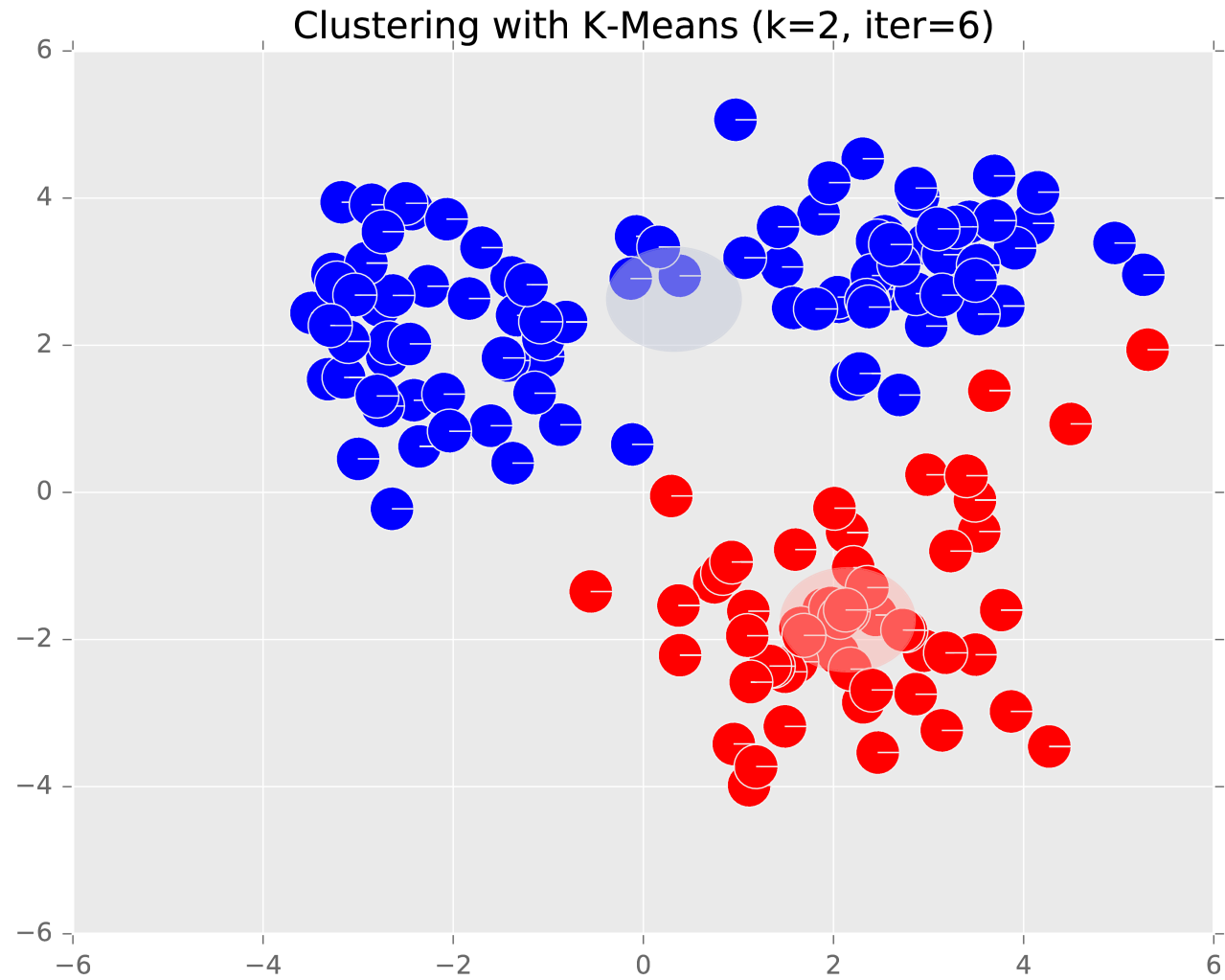
K-means:
Example
($K = 2$)



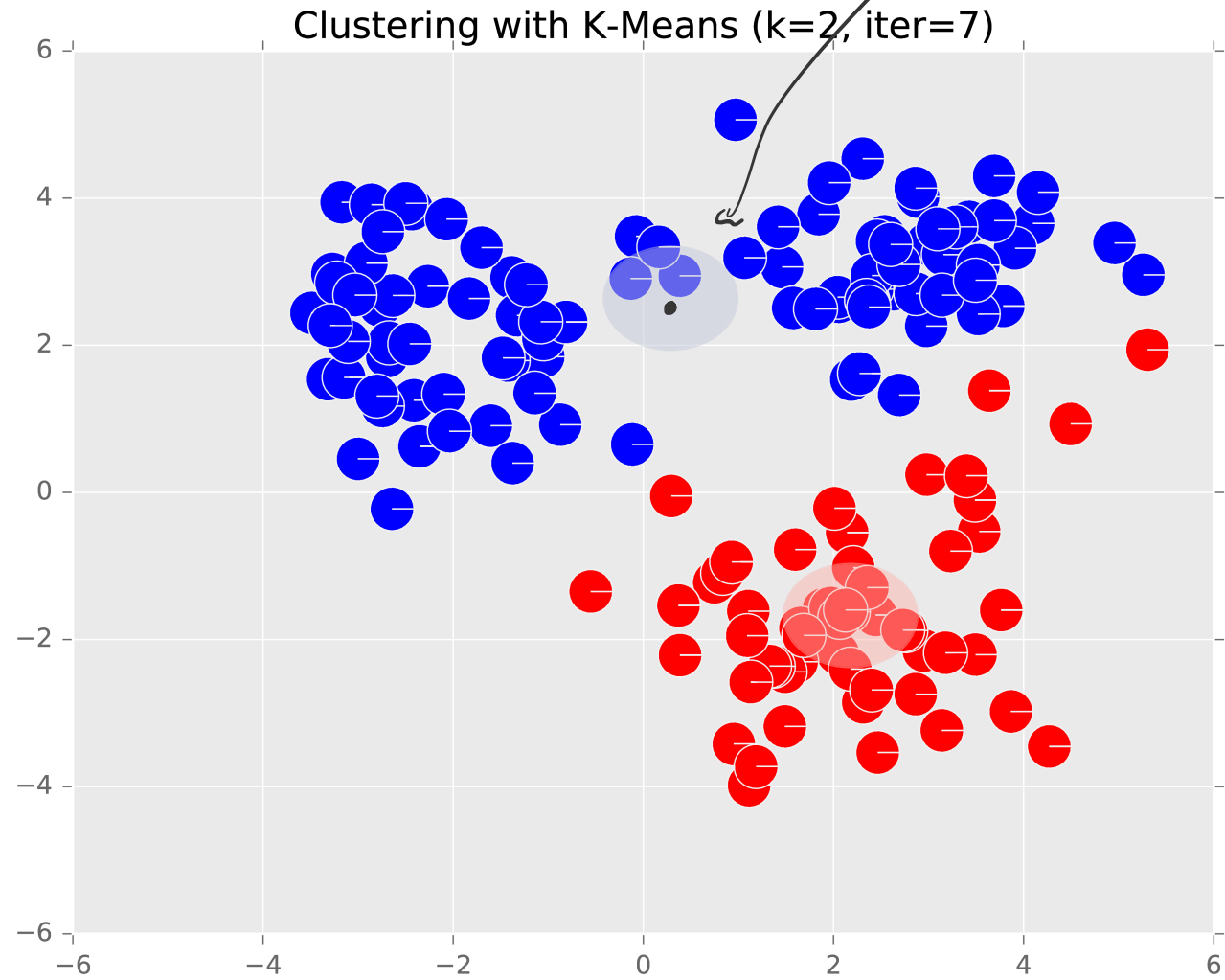
K-means:
Example
($K = 2$)



K-means:
Example
($K = 2$)

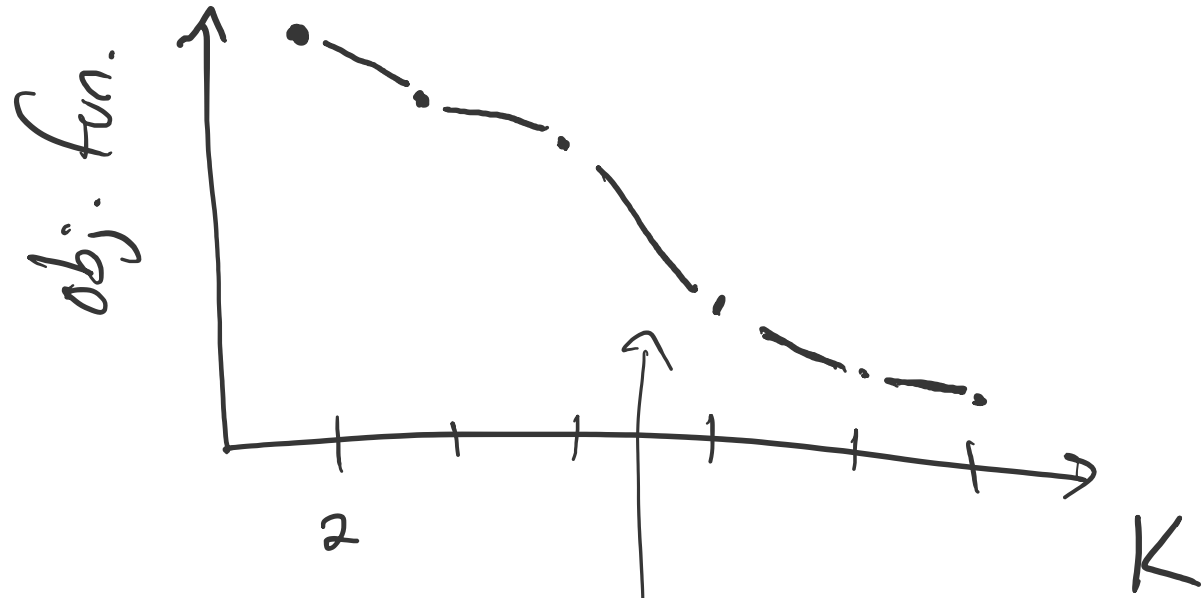


K-means: Example (*K* = 2)



Setting K

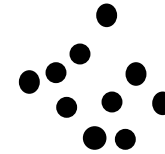
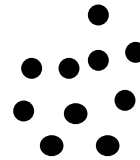
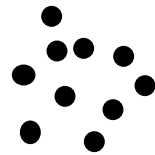
- Idea: choose the value of K that minimizes the objective function



look for big dips in the
obj. fun.

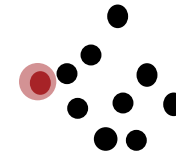
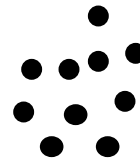
Initializing K -means

- Common choice: choose K data points at random to be the initial cluster centers (Lloyd's method)



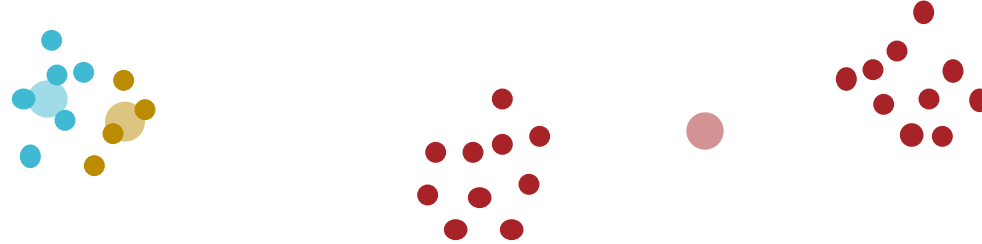
Initializing K -means

- Common choice: choose K data points at random to be the initial cluster centers (Lloyd's method)



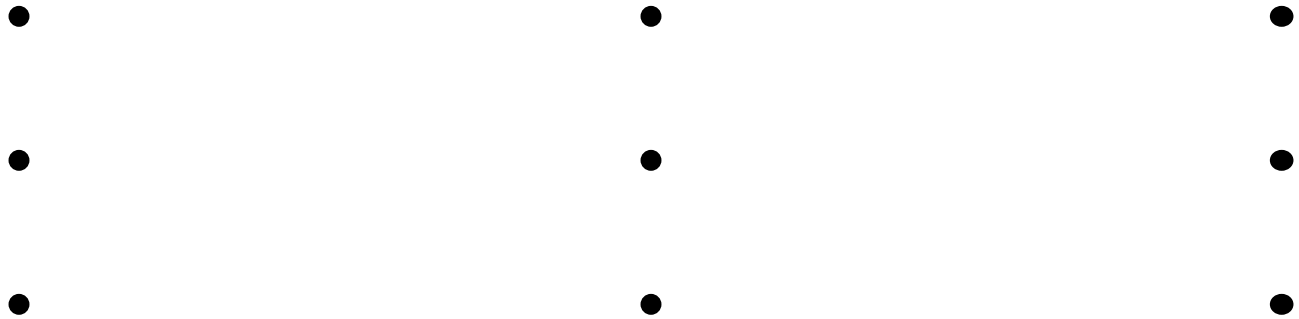
Initializing K -means

- Common choice: choose K data points at random to be the initial cluster centers (Lloyd's method)



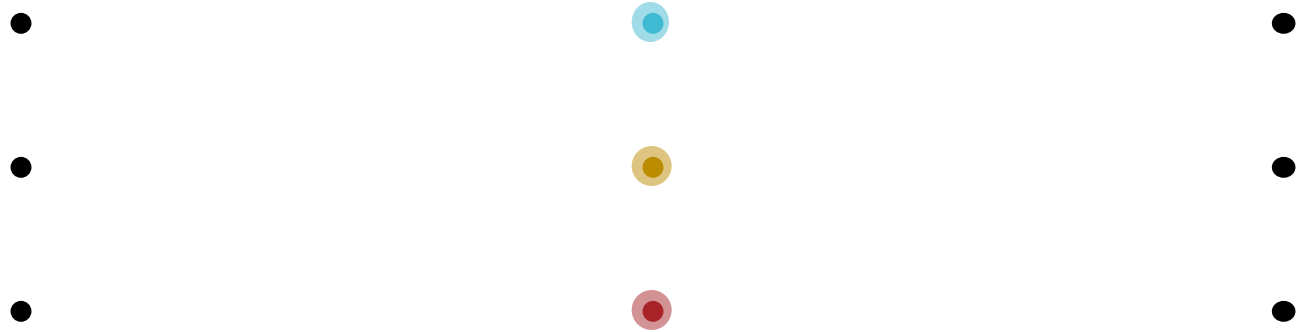
Initializing K -means

- Common choice: choose K data points at random to be the initial cluster centers (Lloyd's method)



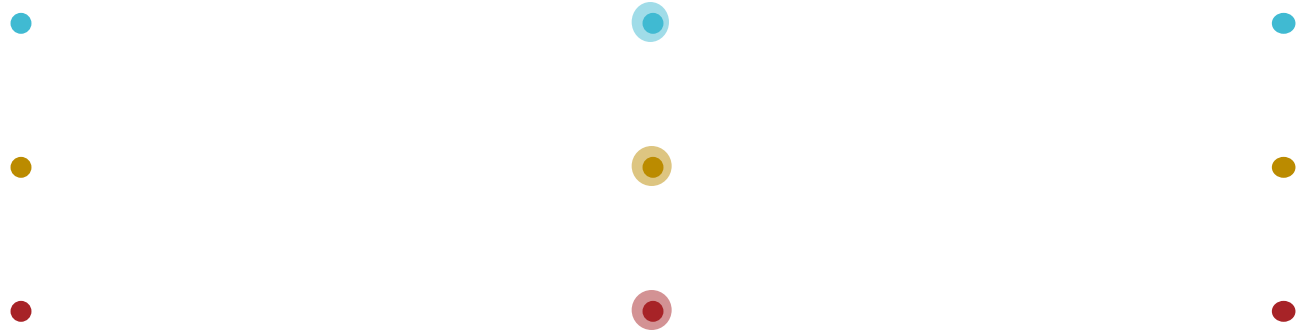
Initializing K -means

- Common choice: choose K data points at random to be the initial cluster centers (Lloyd's method)



Initializing K -means

- Common choice: choose K data points at random to be the initial cluster centers (Lloyd's method)



- Lloyd's method converges to a local minimum and that local minimum can be arbitrarily bad (relative to the optimal clusters)
- Intuition: want initial cluster centers to be far apart from one another

K -means++ (Arthur and Vassilvitskii, 2007)

1. Choose the first cluster center randomly from the data points.
2. For each other data point \mathbf{x} , compute $D(\mathbf{x})$, the distance between \mathbf{x} and the closest cluster center.
3. Select the next cluster center proportional to $D(\mathbf{x})^2$.
4. Repeat 2 and 3 $K - 1$ times.
 - K -means++ achieves a $O(\log K)$ approximation to the optimal clustering in expectation
 - Both Lloyd's method and K -means++ can benefit from multiple random restarts.

Key Takeaways

- K -means objective function & model parameters
- Block-coordinate descent
- Setting K
- Initializing K means