

# 10-301/601: Introduction to Machine Learning

## Lecture 23: Clustering

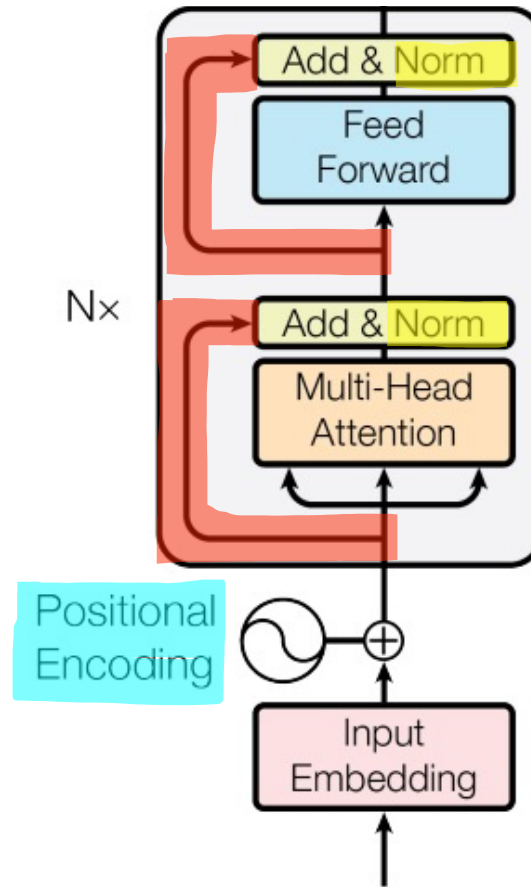
Henry Chai

6/4/25

# Front Matter

- Announcements
  - HW5 released on 6/3, due 6/6 at 11:59 PM
  - Schedule change: two recitations this week
    - **Recitation on 6/4 (today!) will be a PyTorch tutorial**
    - Recitation on 6/5 will be Quiz 3 preparation

# Recall: Transformers



- In addition to multi-head attention, transformer architectures use
  1. Positional encodings
  2. Layer normalization
  3. Residual connections
  4. A fully-connected feed-forward network

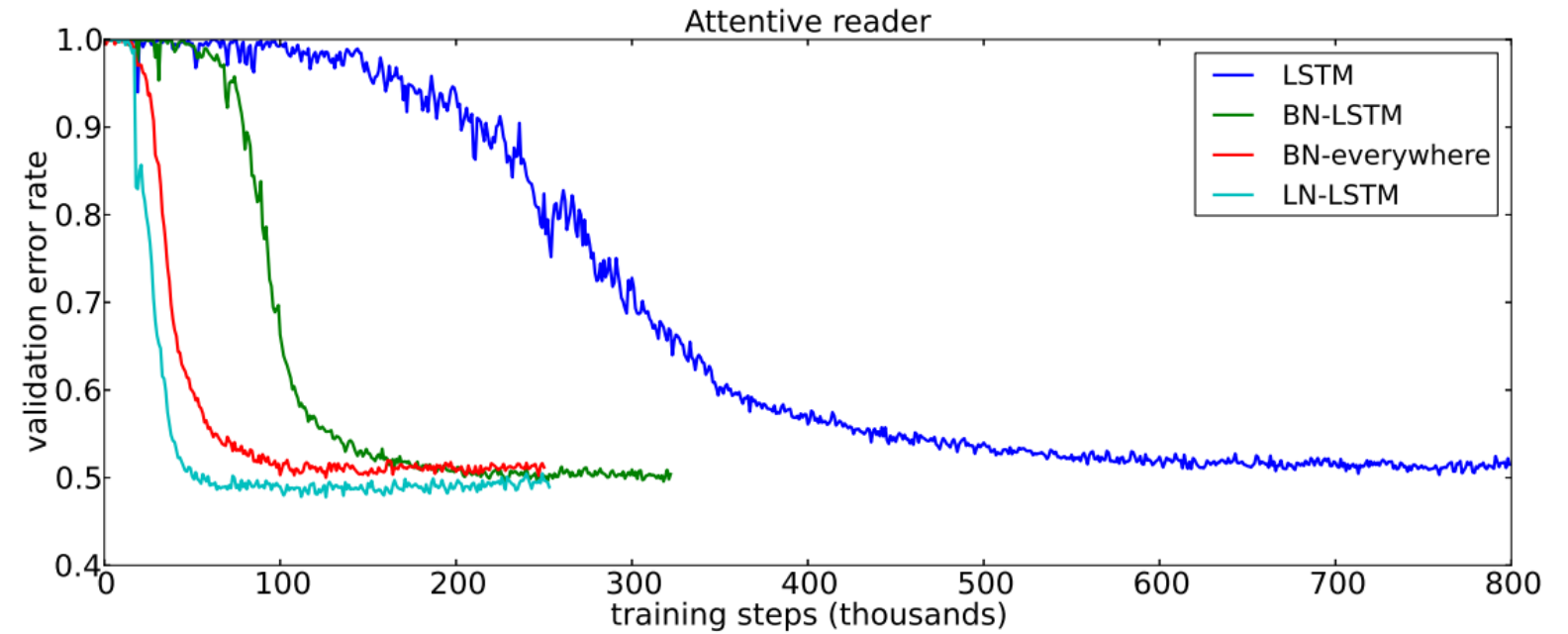
# Layer Normalization

- Issue: for certain activation functions, the weights in later layers are **highly sensitive** to changes in the earlier layers
  - Small changes to weights in early layers are amplified so weights in deeper layers have to deal with massive dynamic ranges → slow optimization convergence
- Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$

$$H' = \gamma \left( \underbrace{\frac{H - \mu}{\sigma}} \right) + \beta$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the values in the vector  $H$

# Layer Normalization



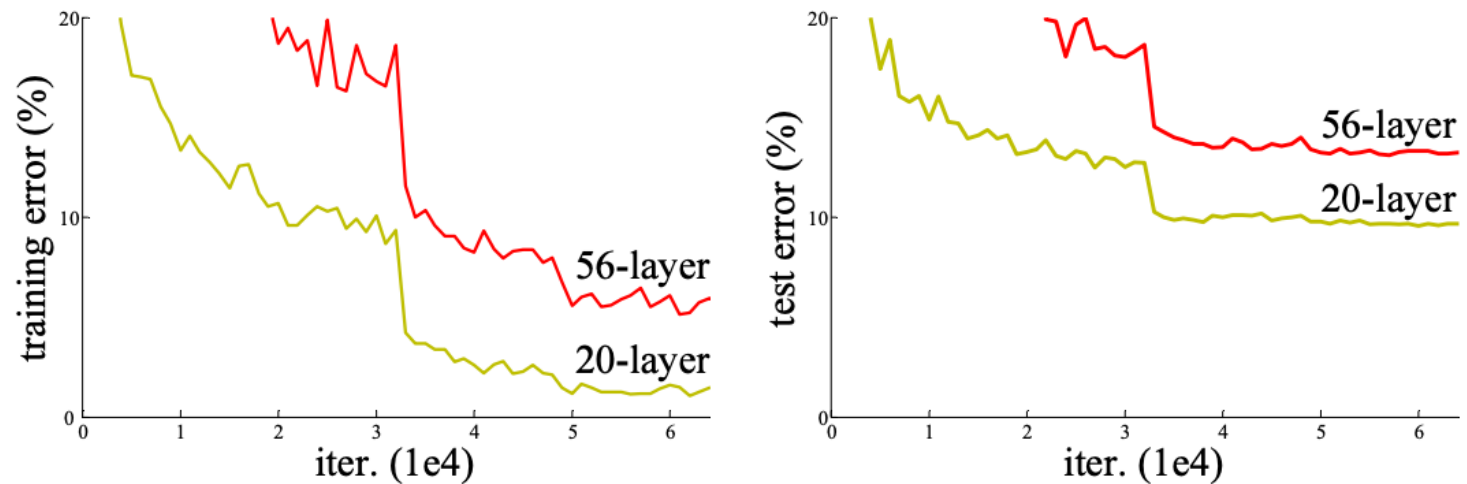
- Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$

$$H' = \gamma \left( \frac{H - \mu}{\sigma} \right) + \beta$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the values in the vector  $H$

# Residual Connections

- Observation: early deep neural networks suffered from the “degradation” problem where adding more layers actually made performance worse!



- Wait but this is ridiculous: if the later layers aren't helping, couldn't they just learn the identity transformation???
- Insight: neural network layers actually have a hard time learning the identity function

# Residual Connections

- Observation: early deep neural networks suffered from the “degradation” problem where adding more layers actually made performance worse!
- Idea: add the input embedding back to the output of a layer

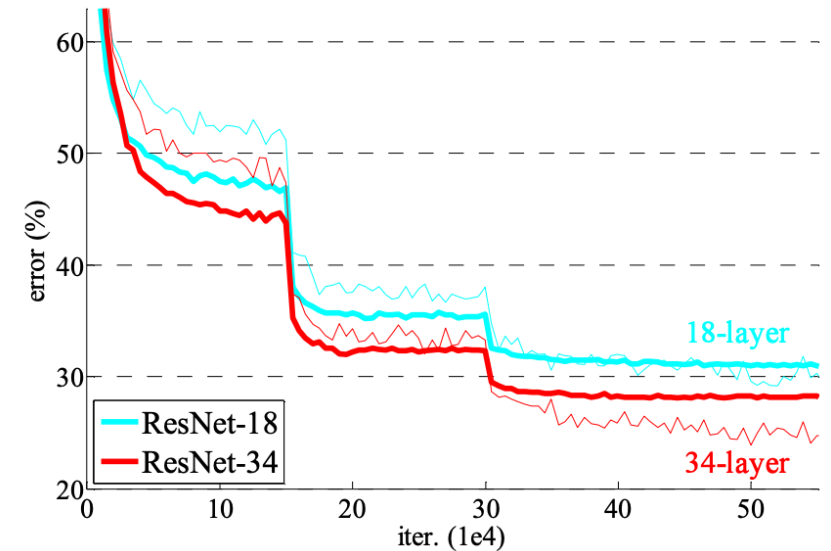
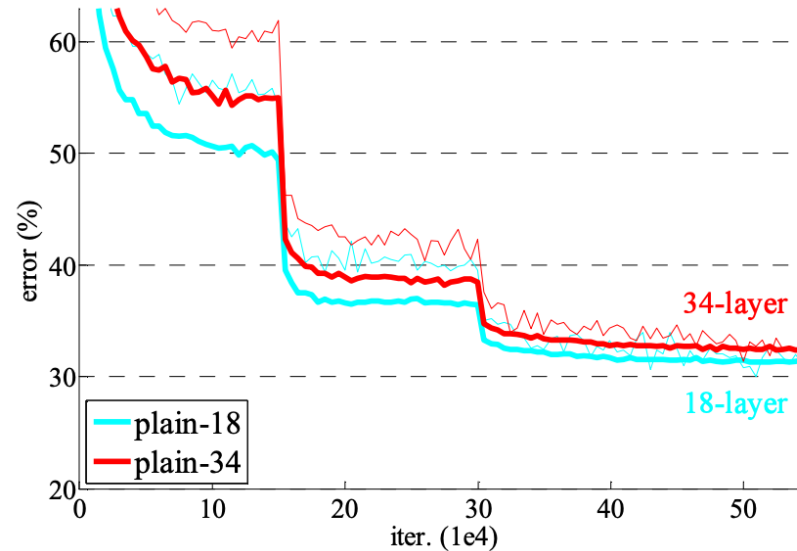
$$H' = H(x^{(i)}) + x^{(i)}$$

- Suppose the target function is  $f$ 
  - Now instead of having to learn  $f(x^{(i)})$ , the hidden layer just needs to learn the residual  $r = f(x^{(i)}) - x^{(i)}$
  - If  $f$  is the identity function, then the hidden layer just needs to learn  $r = 0$ , which is easy for a neural network!

# Residual Connections

- Observation: early deep neural networks suffered from the “degradation” problem where adding more layers actually made performance worse!
- Idea: add the input embedding back to the output of a layer

$$H' = H(x^{(i)}) + x^{(i)}$$





# 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

# 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:
  - clustering of syllable pronunciations
  - clustering novels/texts
  - clustering base pairs in genes

## 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 *and why*?

0 surveys completed



0 surveys underway

## Which partition or clustering do you prefer?

Option A

Option B



## Justify your response to the previous question

Join by Web

**[Pollev.com/301601polls](https://Pollev.com/301601polls)**

Join by QR code

Scan with your camera app



# General Recipe for Machine Learning

- Define a model and model parameters

– Assume  $K$  clusters  $\hookrightarrow$  similarity is defined by Euclidean distance

– Parameters: Assignments –  $z^{(1)}, z^{(2)}, \dots, z^{(N)}$

- Write down an objective function

centers –  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_K$

$$\text{minimize } J(z^{(1)}, \dots, z^{(N)}, \vec{\mu}_1, \dots, \vec{\mu}_K) = \sum_{n=1}^N \|\vec{x}^{(n)} - \vec{\mu}_{z^{(n)}}\|_2$$

- Optimize the objective w.r.t. the model parameters

– Block coordinate descent

# Recipe for $K$ -means

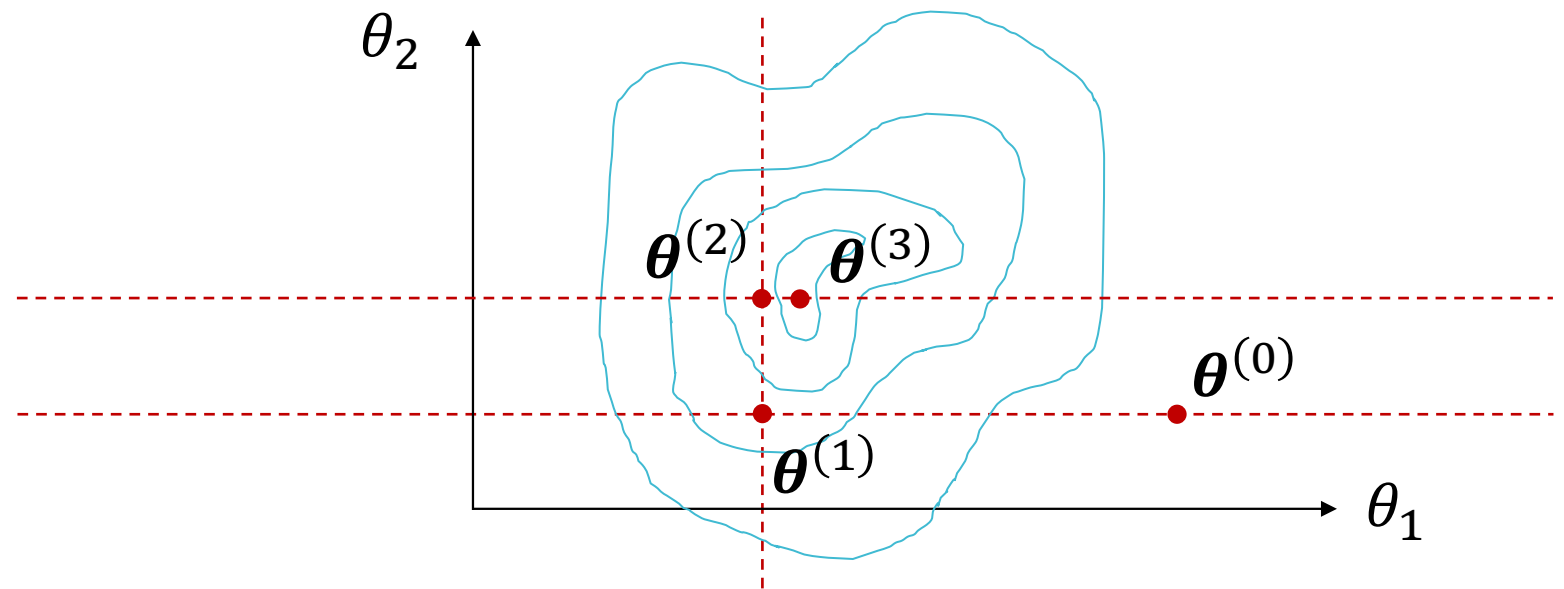
- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

# 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 ( $\alpha$  or  $\beta$ ) 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

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

- If  $\mu_1, \dots, \mu_K$  are fixed

$$\hat{z}^{(n)} = \operatorname{argmin}_{c \in \{1, \dots, K\}} \|\vec{x}^{(n)} - \vec{\mu}_c\|_2$$

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

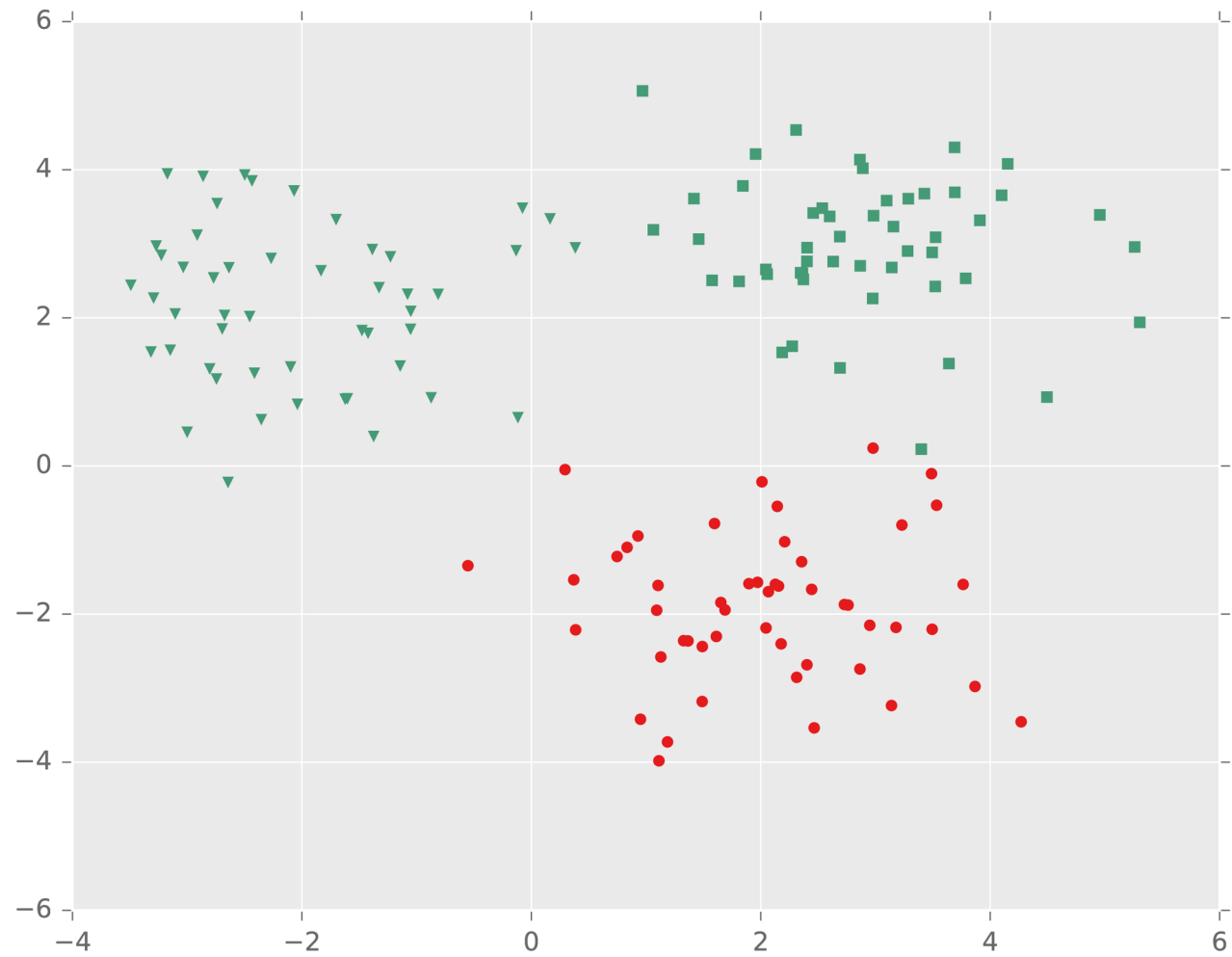
$$\hat{\vec{\mu}}_c = \frac{1}{N_c} \sum_{n: z^{(n)} = c} \vec{x}^{(n)}$$

where  $N_c$  is the # of data points assigned to cluster  $C$

# K-means Algorithm

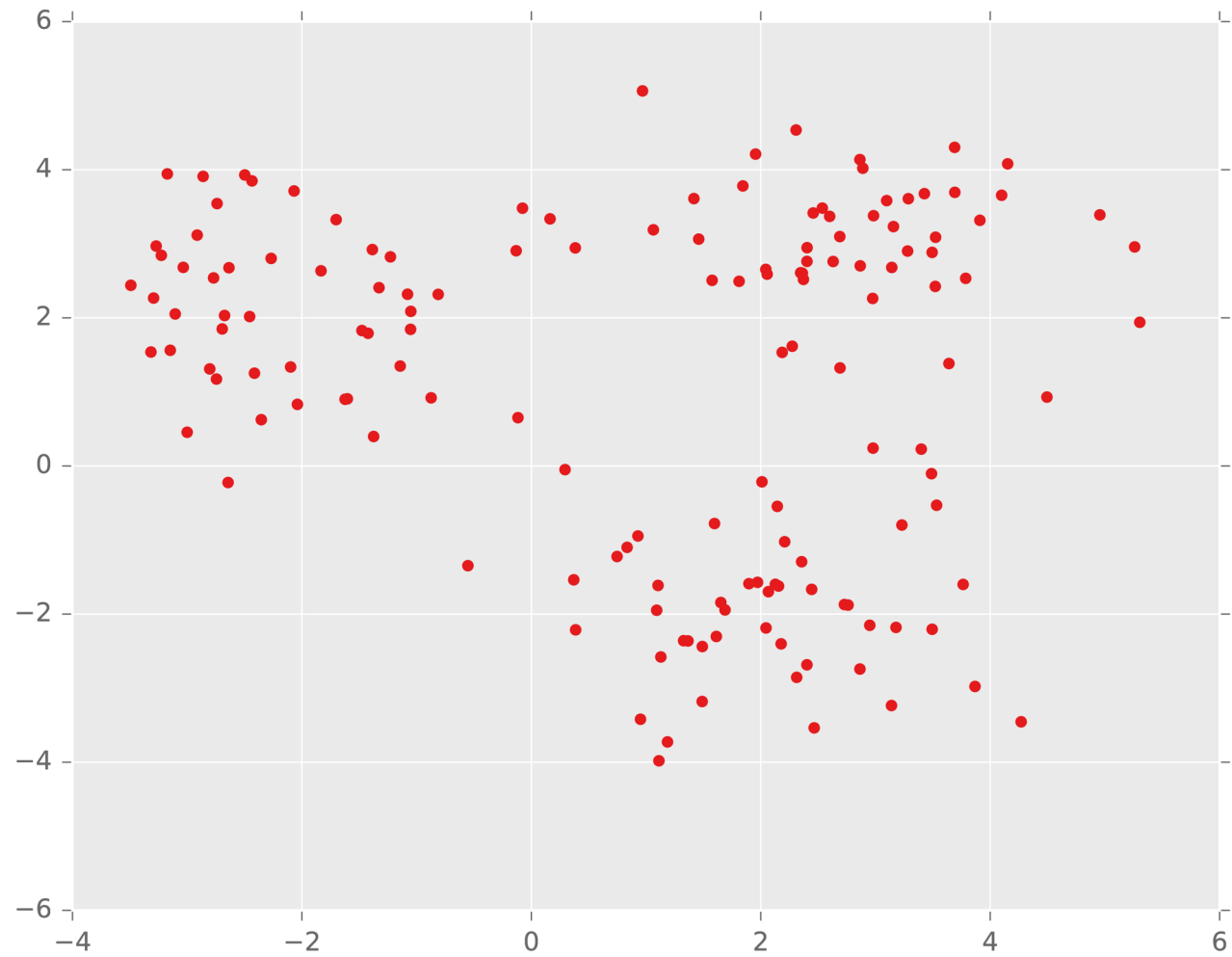
- Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, K$ 
  1. Initialize cluster centers  $\mu_1, \dots, \mu_K$
  2. While NOT CONVERGED
    - a. Assign each data point to the cluster with the nearest cluster center:
$$z^{(n)} = \underset{k}{\operatorname{argmin}} \|\mathbf{x}^{(n)} - \mu_k\|_2$$
    - b. Recompute the cluster centers:
$$\mu_k = \frac{1}{N_k} \sum_{n: z^{(n)}=k} \mathbf{x}^{(n)}$$
where  $N_k$  is the number of data points in cluster  $k$
- Output: cluster centers  $\mu_1, \dots, \mu_K$  and cluster assignments  $z^{(1)}, \dots, z^{(N)}$

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

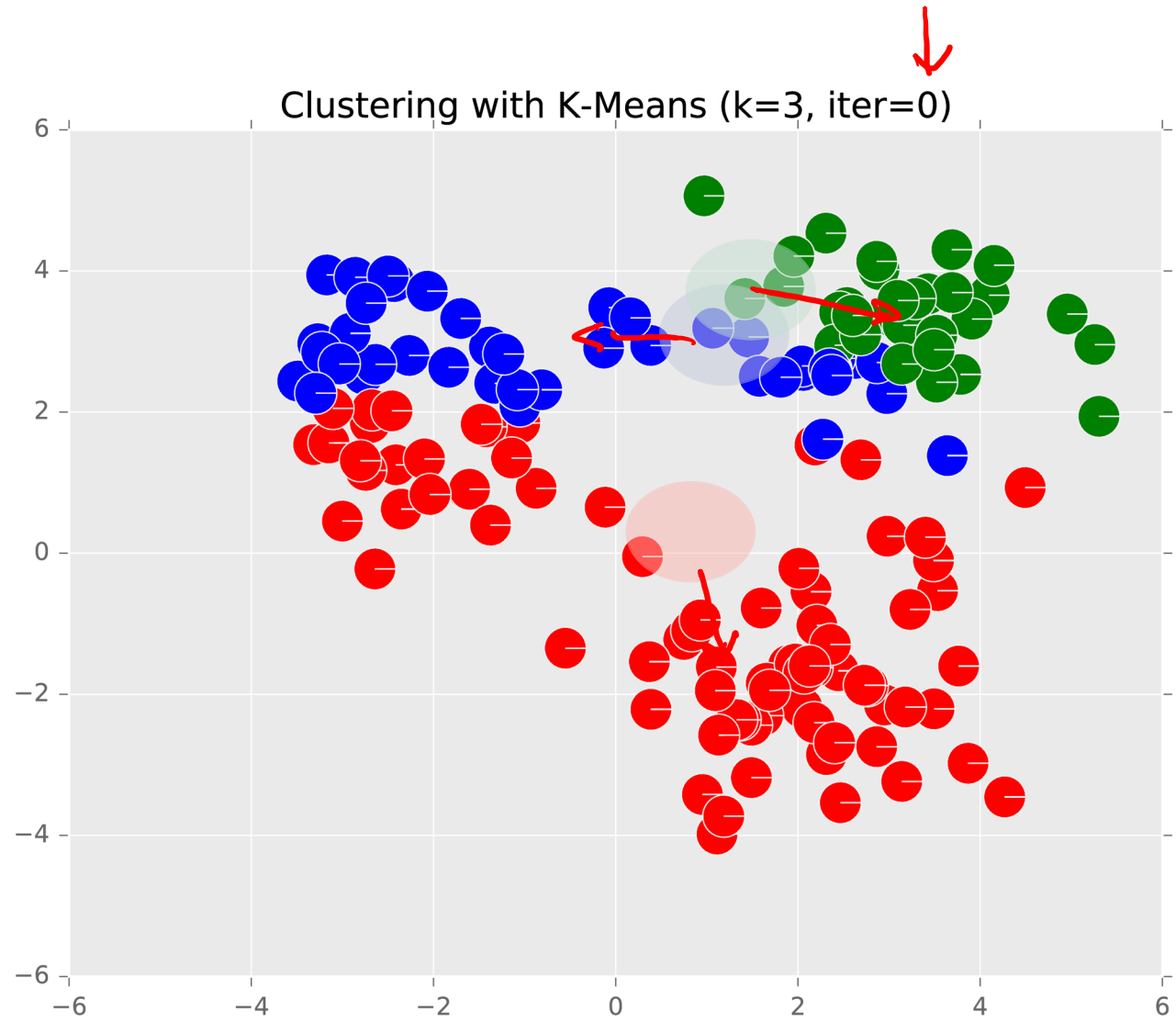




# $K$ -means: Example ( $K = 3$ )



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



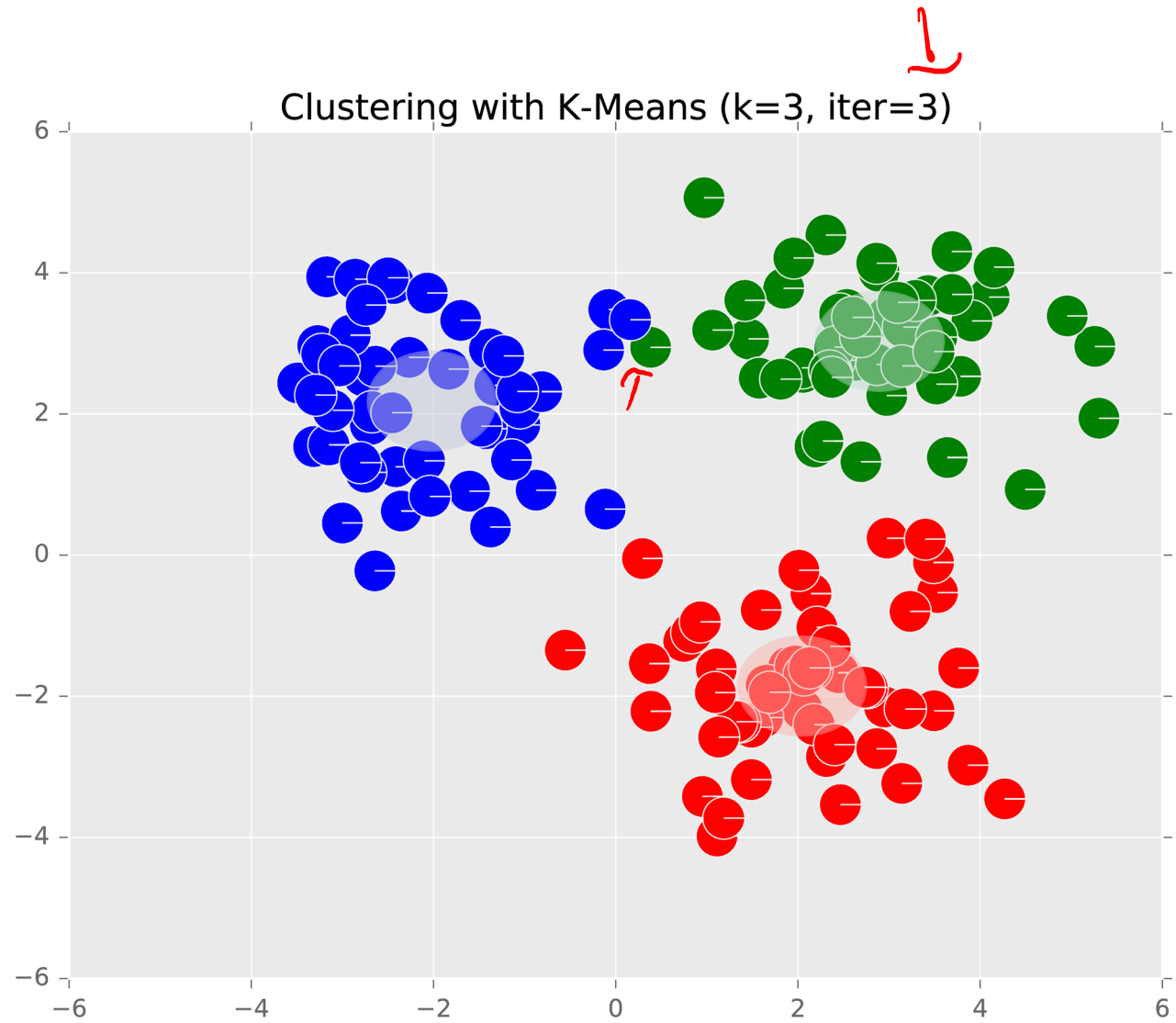
# $K$ -means: Example ( $K = 3$ )



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



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



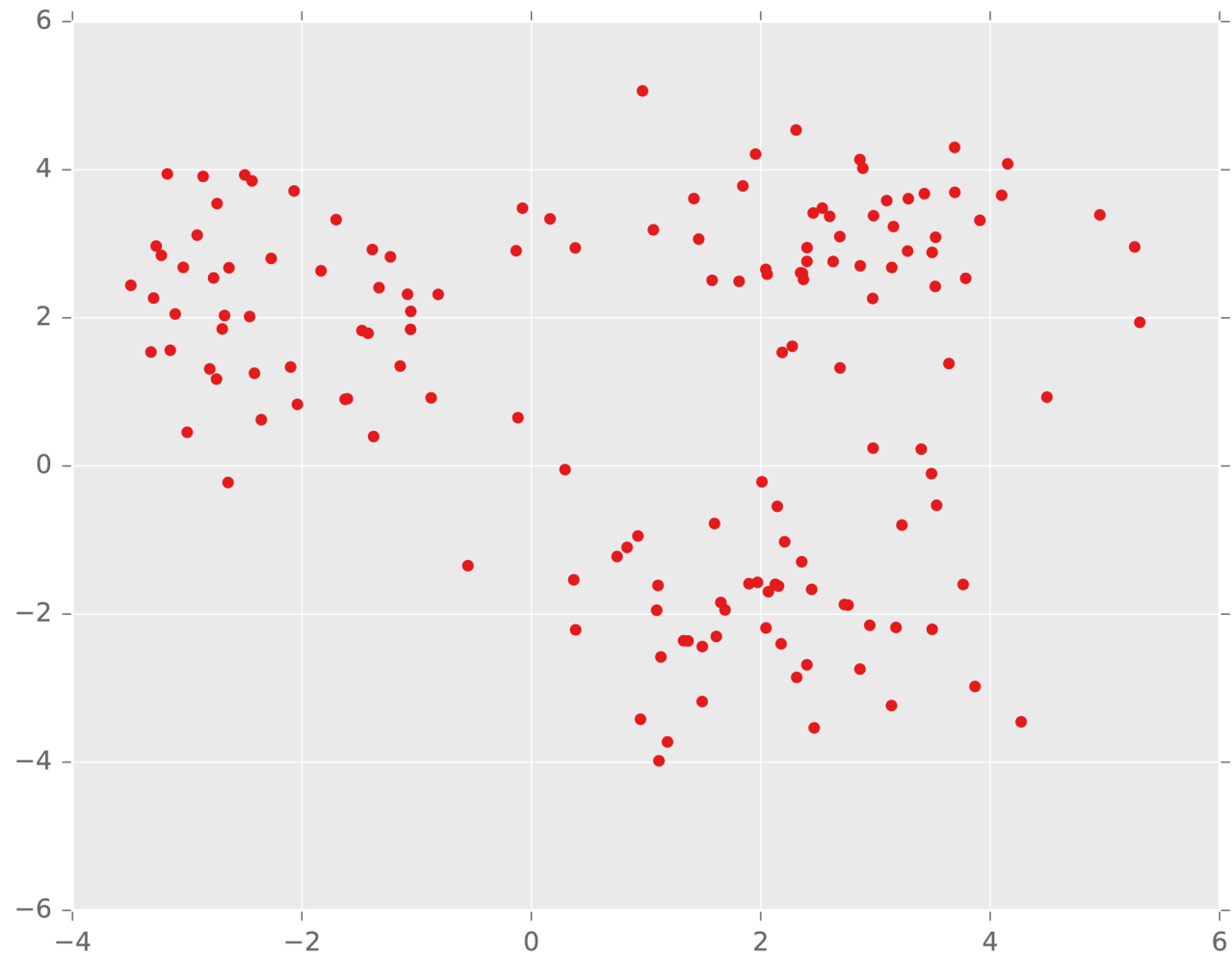
# $K$ -means: Example ( $K = 3$ )



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

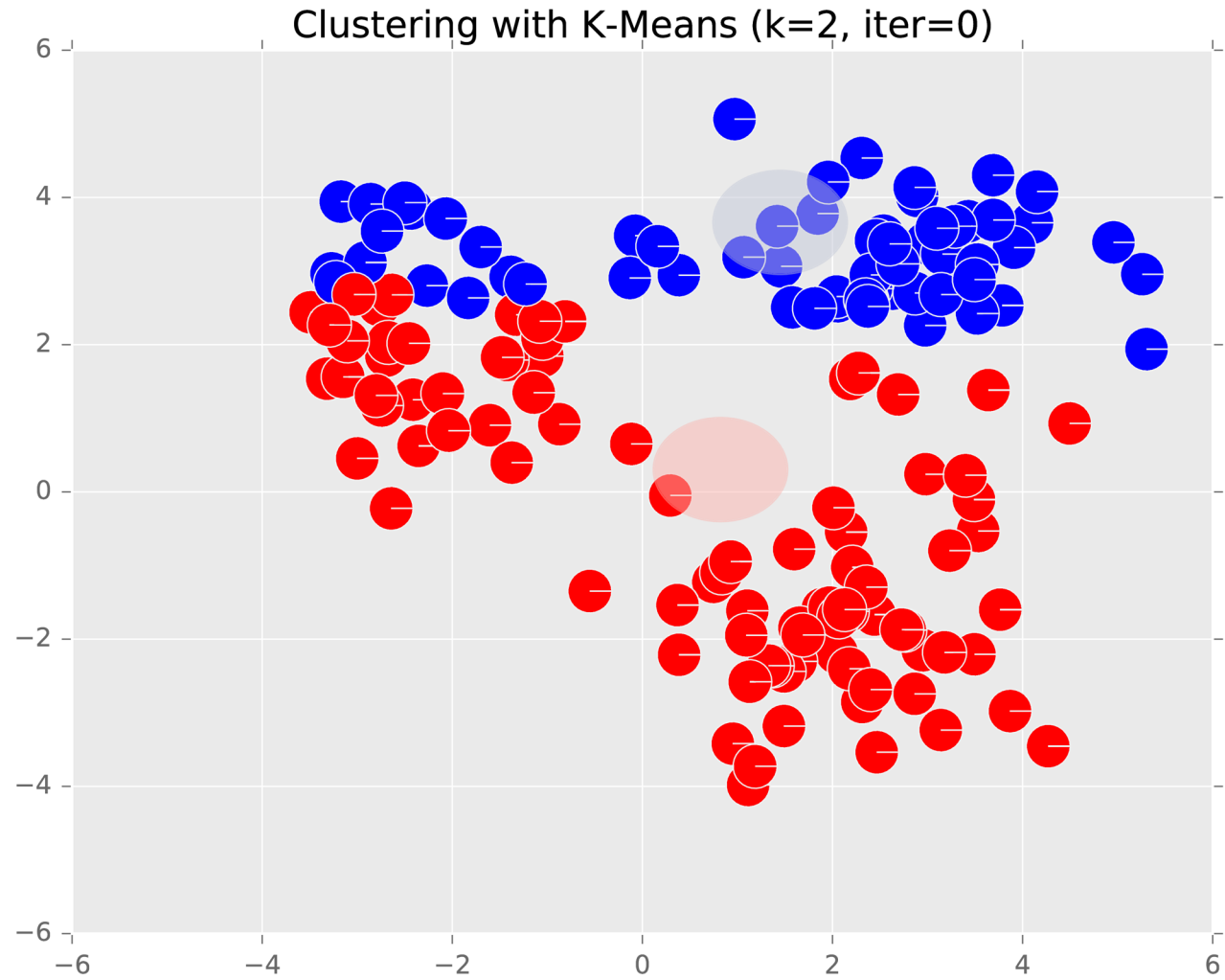


# $K$ -means: Example ( $K = 2$ )

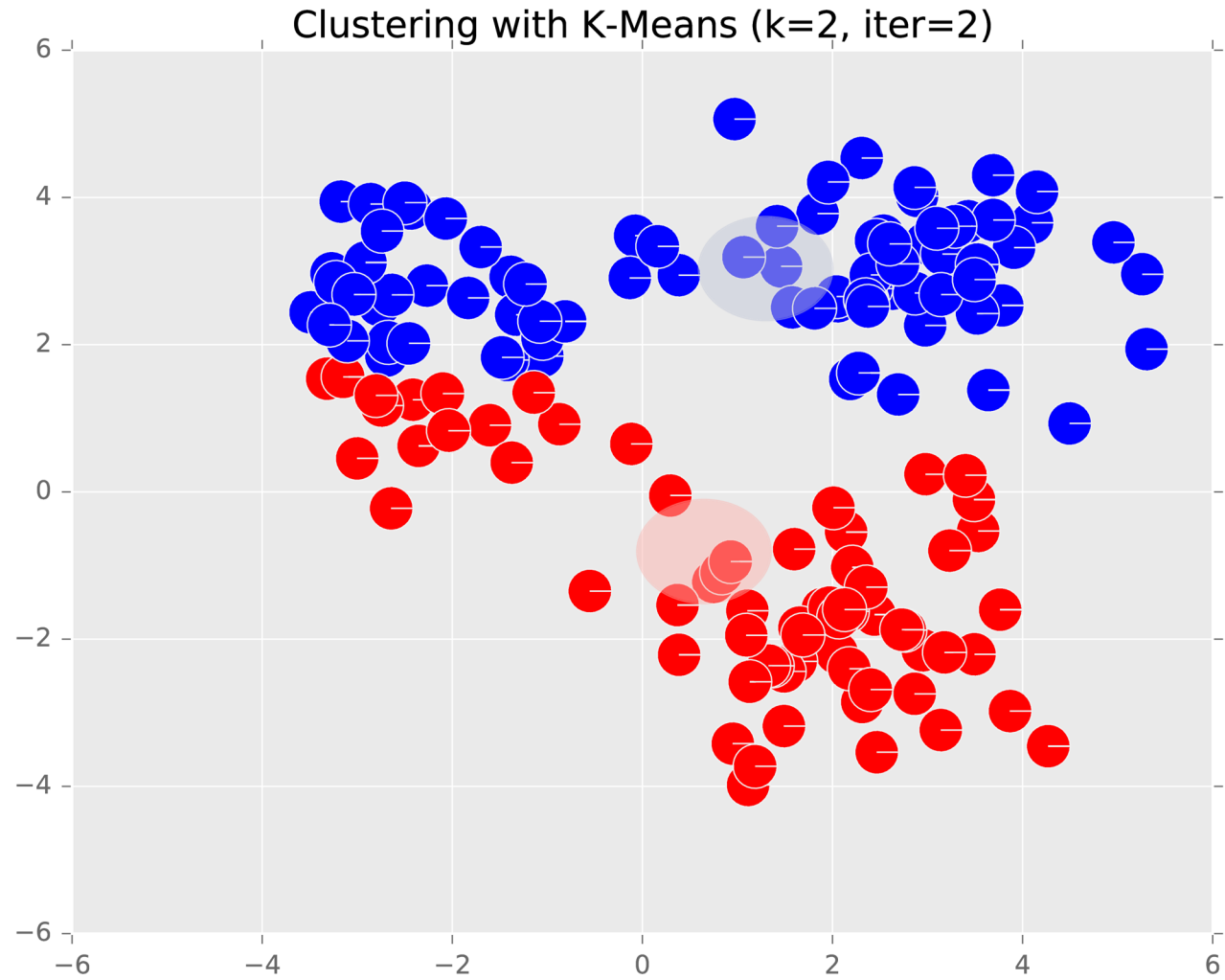




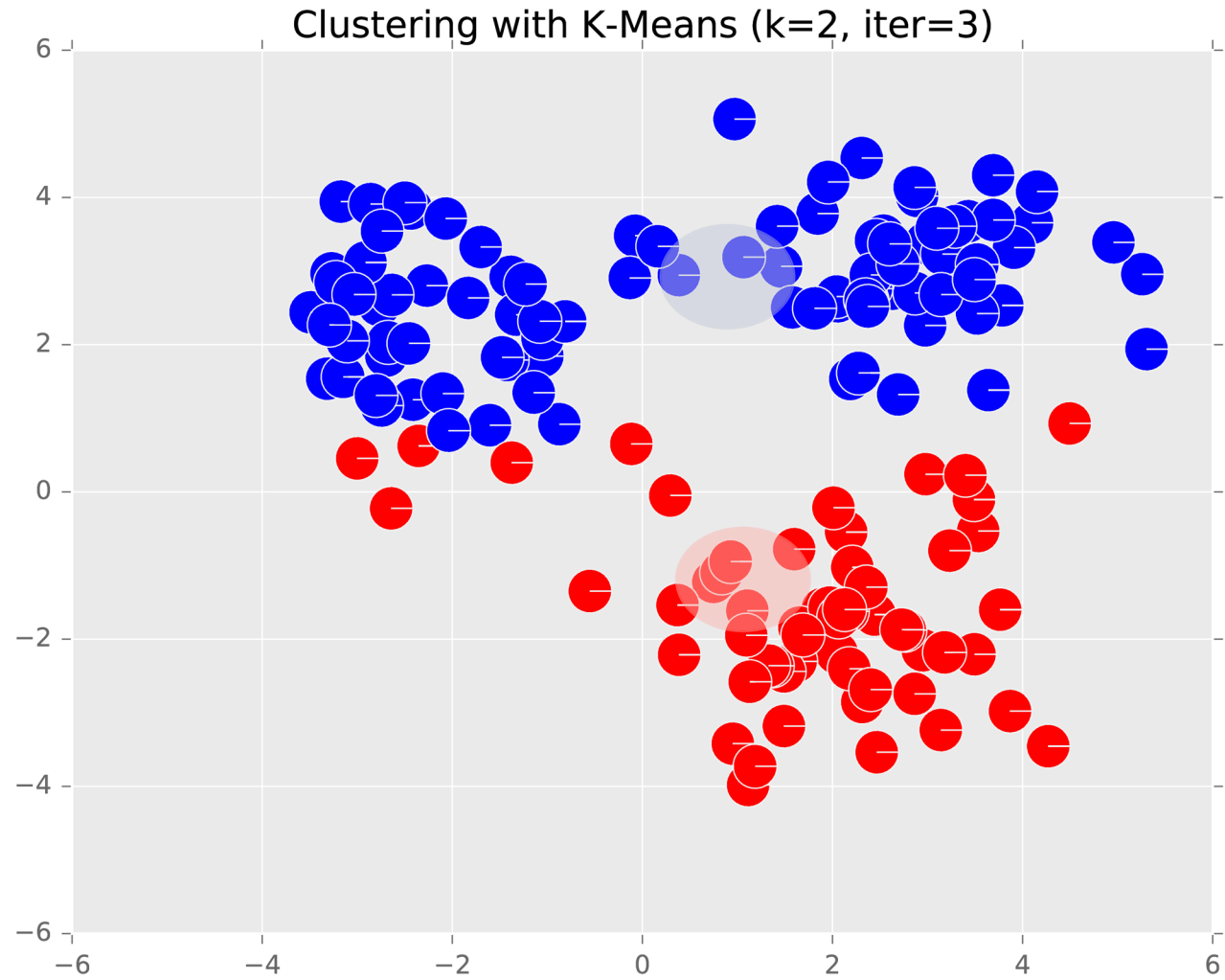
# $K$ -means: Example ( $K = 2$ )



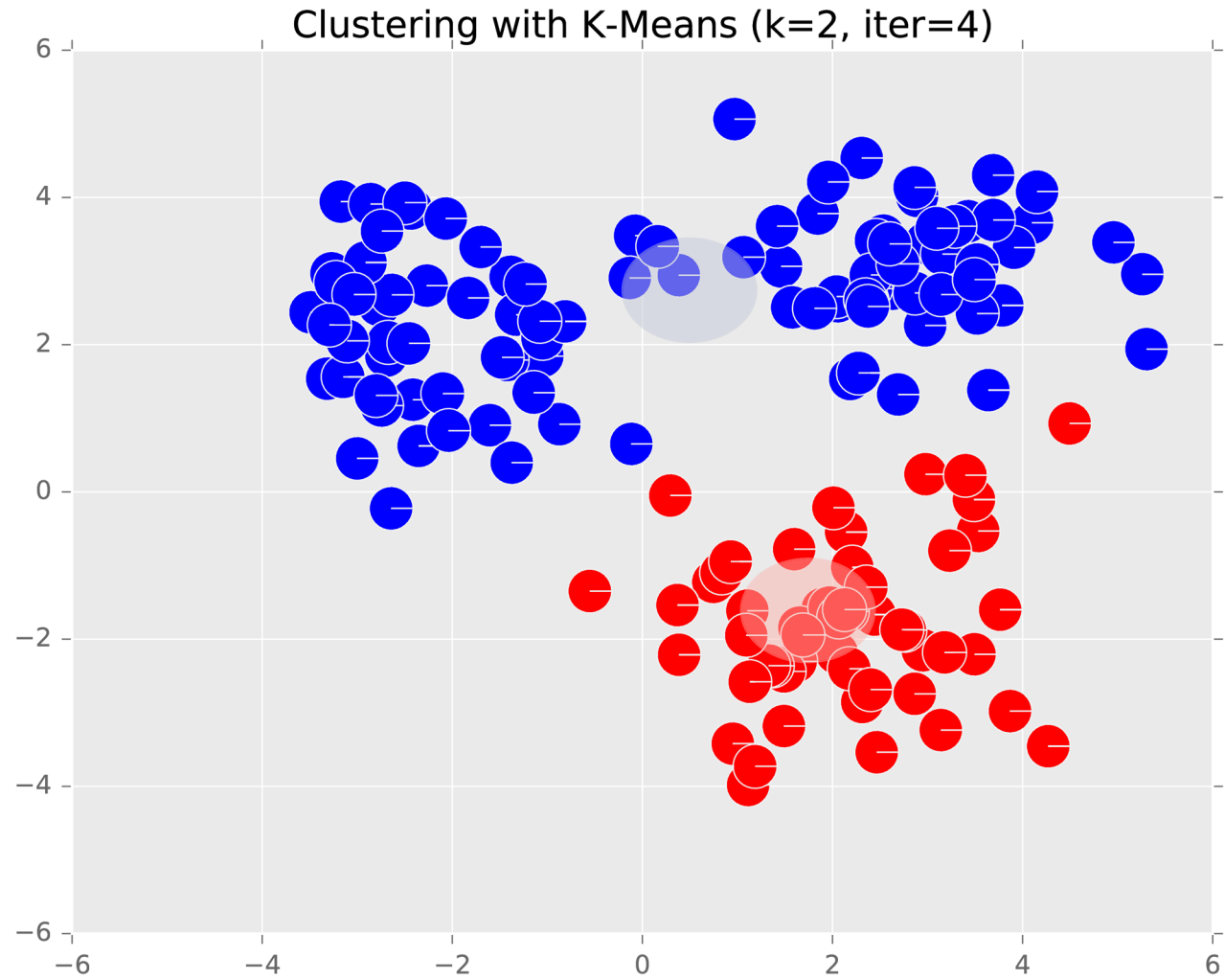
# $K$ -means: Example ( $K = 2$ )



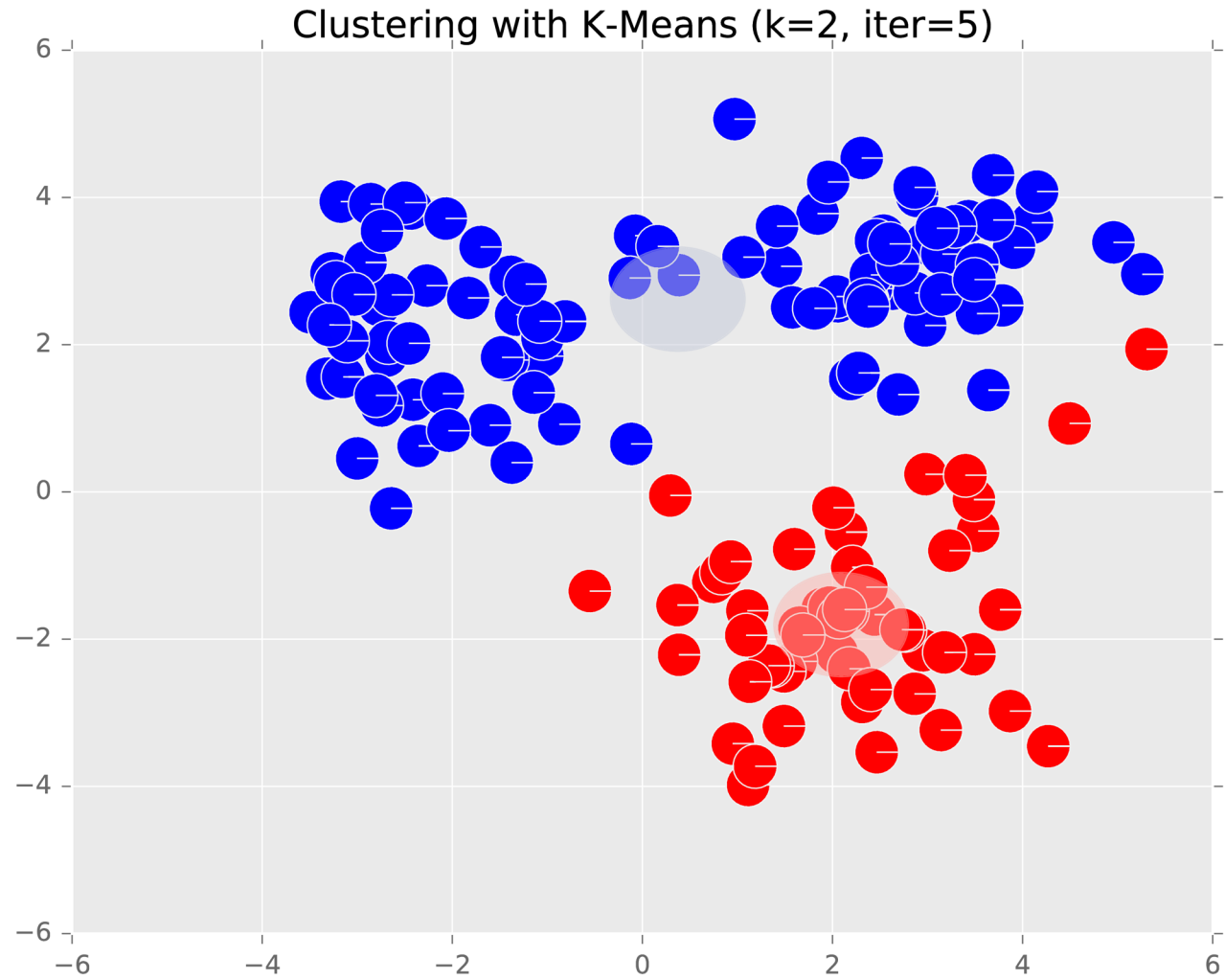
# $K$ -means: Example ( $K = 2$ )



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



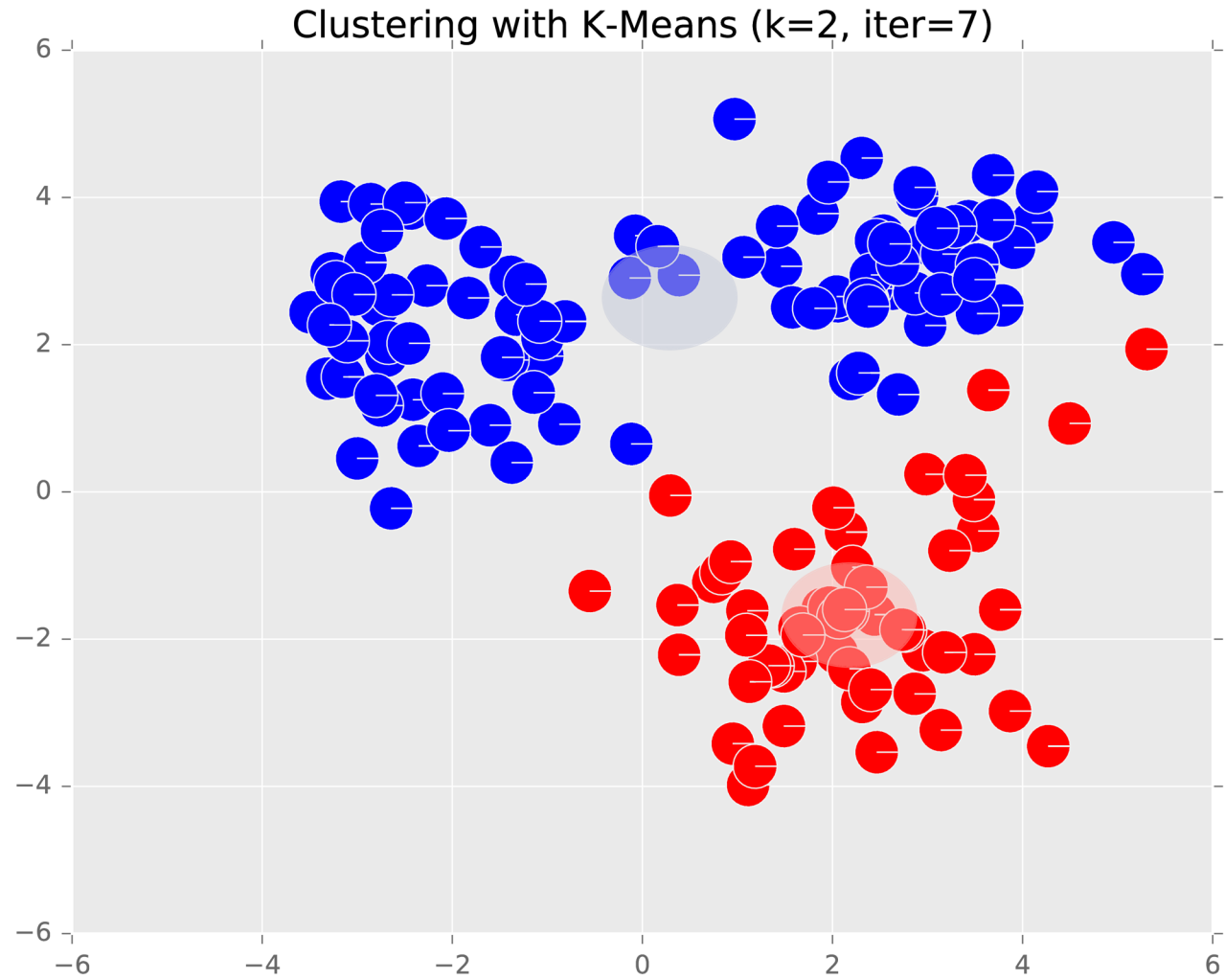
# $K$ -means: Example ( $K = 2$ )



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



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



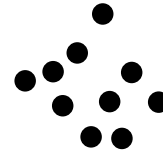
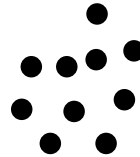
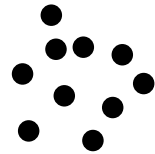
## Setting $K$

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



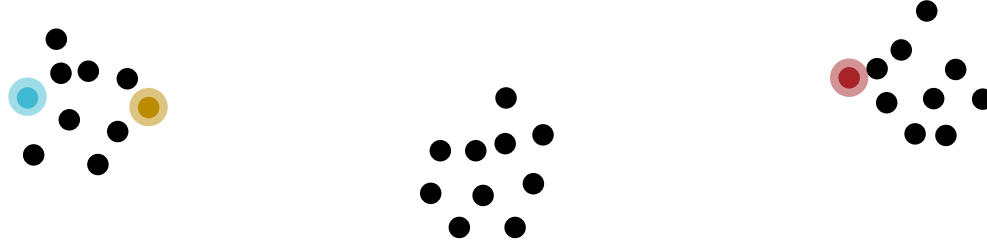
# Initializing $K$ -means

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



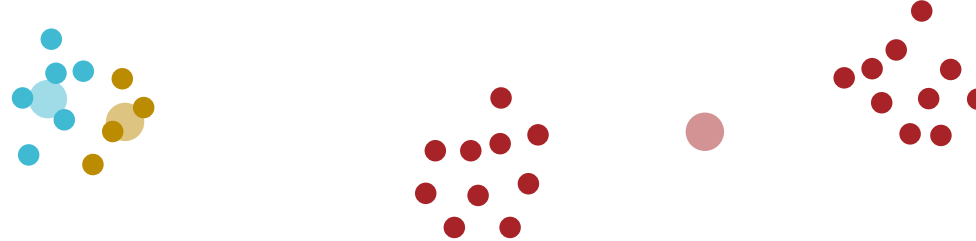
# Initializing $K$ -means

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



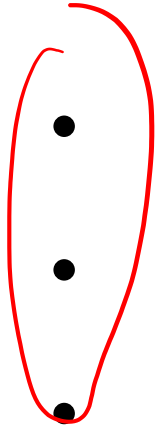
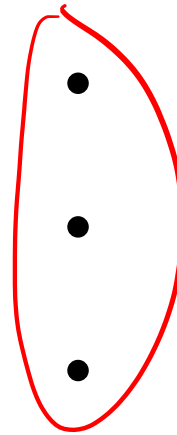
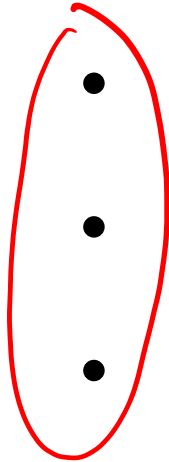
# Initializing $K$ -means

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



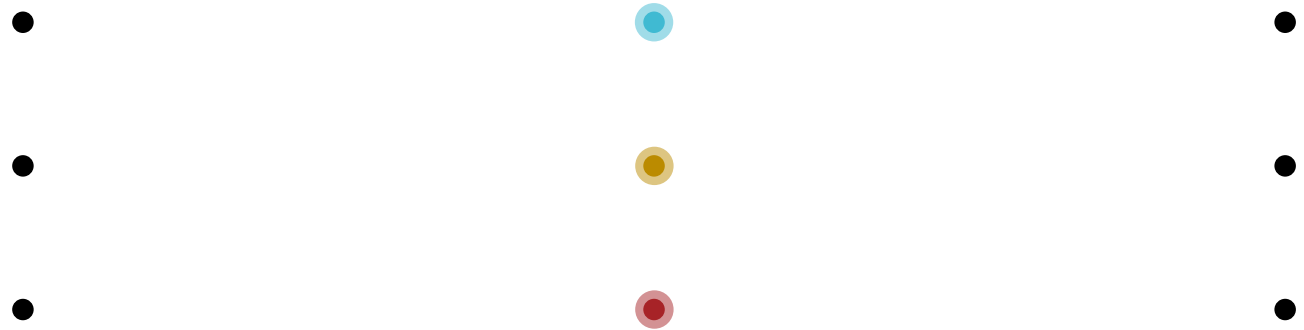
# Initializing $K$ -means

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



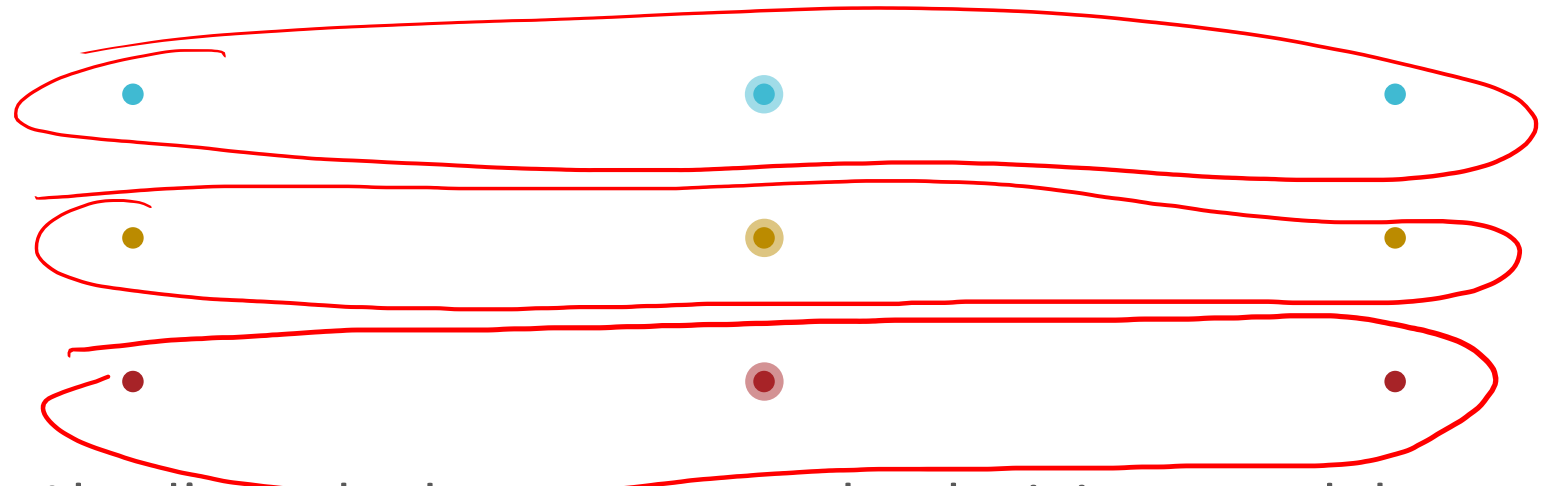
# Initializing $K$ -means

- 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