

10-701: Introduction to Machine Learning

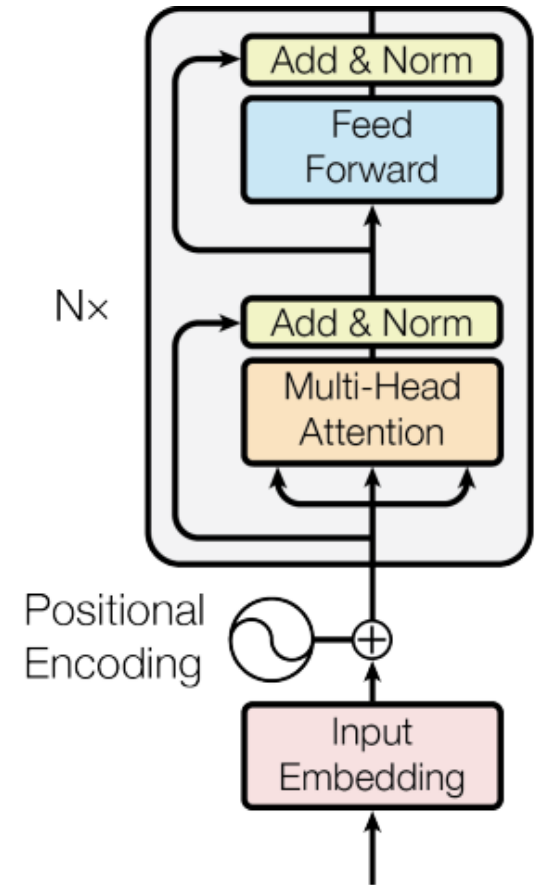
Lecture 14 – Unsupervised Learning: Clustering

Hoda Heidari

* Slides adopted from F24 offering of 10701 by Henry Chai.

Transformers

1. **Embed words:** Words are turned into vectors (points in a space that roughly encode semantic meaning).
2. **Attention Scoring:** Each word looks at all other words and asks: “How relevant are you to me given the task at hand?”
3. **Blend information:** Each word builds a new representation by taking a weighted mix of the other words, weighted by their scores.
4. **Repeat:** Stack that many times-- layers repeat the same pattern, so deeper layers see richer relational structure/representations.



The Attention Mechanism

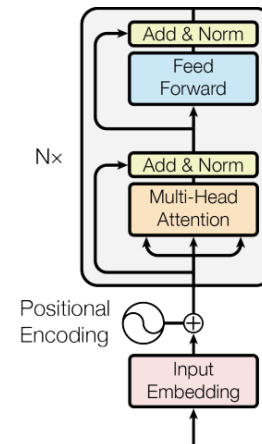
Analogy: attention as soft lookup in a dictionary

1. The token x' broadcasts its **query** $q = w_Q^T x'$
2. Every other token x_t offers up its **key** $k_t = W_K x_t$
3. The model computes a *similarity* between the query and keys--scoring how relevant each token is to the query.

$$\underline{s_t(x', x_t)} = \frac{k_t^T q}{\sqrt{\text{length}(q)}}$$

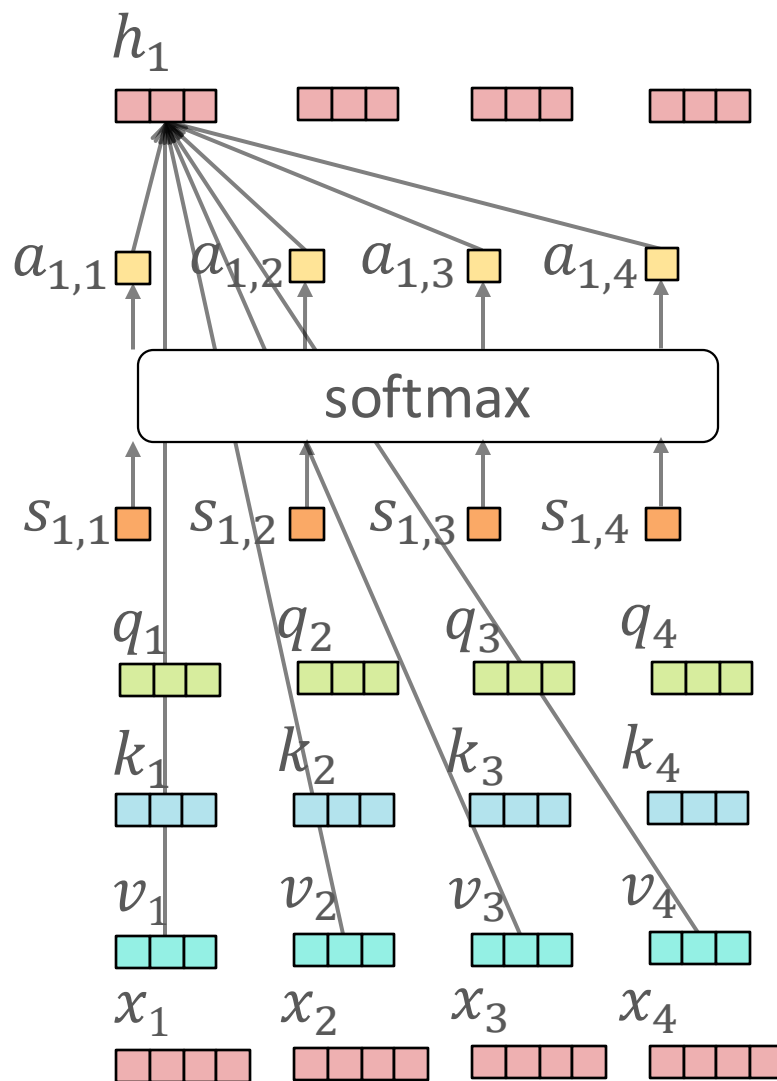
4. Those *scores* are turned into weights: $\text{softmax}(s(x', x_t))$
5. The model returns a weighted sum of the **values** as the contextualized representation:

new representation for x' = $\sum_t \underline{\text{softmax}(s(x', x_t))} \underline{v(x_t)}$



Attention Head

- Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$h_1 = \sum_{j=1}^4 \text{softmax}(s_{1,j}) v_j$$

attention weights

scores: $s_{1,j} = \frac{k_j^T q_1}{\sqrt{\text{length}(k_j)}}$

queries: $q_t = W_Q x_t$

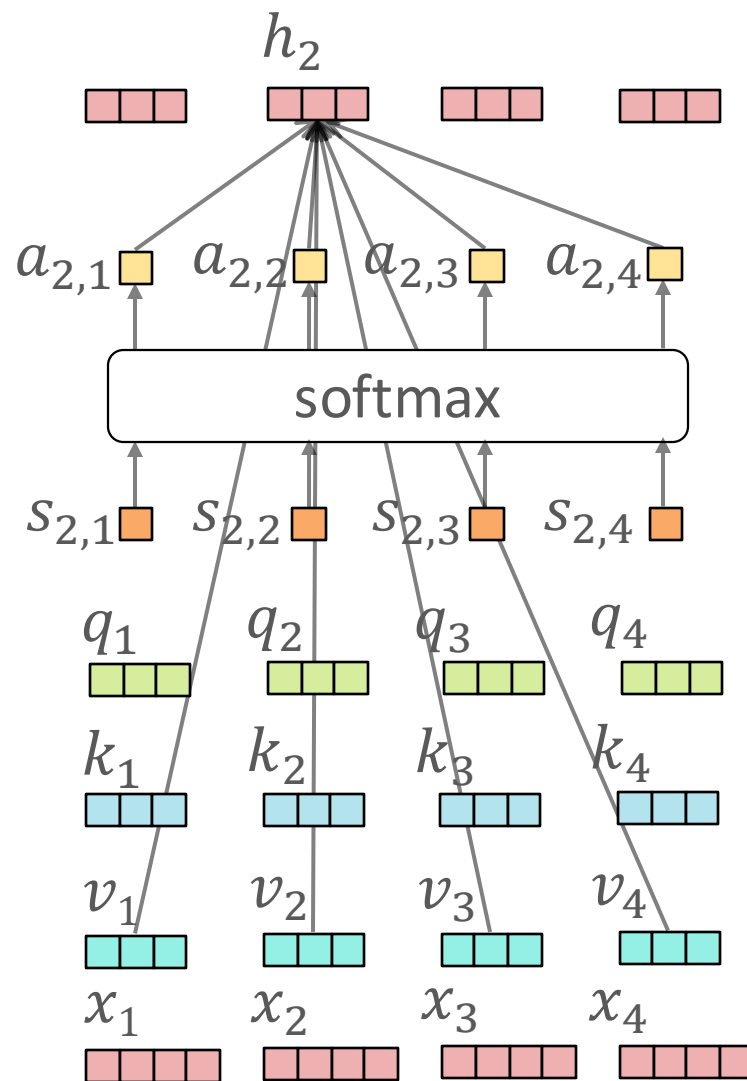
keys: $k_t = W_K x_t$

values: $v_t = W_V x_t$

input tokens

Attention Head

- Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$h_2 = \sum_{j=1}^4 \text{softmax}(s_{2,j}) v_j$$

attention weights

$$\text{scores: } s_{2,j} = \frac{k_j^T q_2}{\sqrt{\text{length}(k_j)}}$$

$$\text{queries: } q_t = W_Q x_t$$

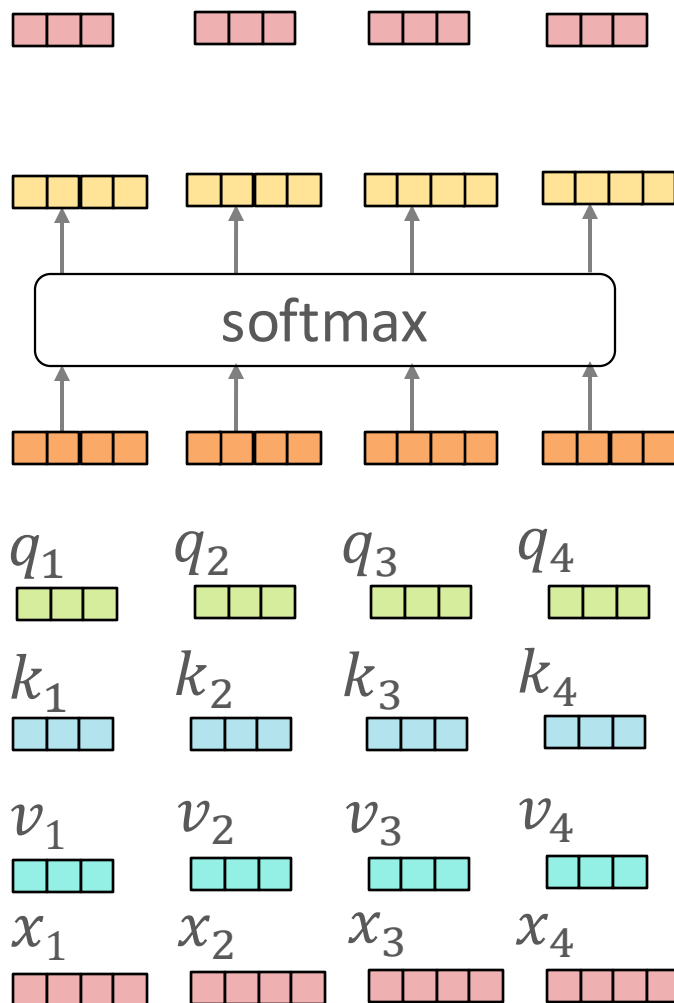
$$\text{keys: } k_t = W_K x_t$$

$$\text{values: } v_t = W_V x_t$$

input tokens

Attention Head: Matrix Form

- Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$H = \text{softmax}(S)V \in \mathbb{R}^{N \times d_v}$$

attention weights

scores: $S = \frac{QK^T}{\sqrt{d_k}} \in \mathbb{R}^{N \times N}$

Handwritten notes: $N \times d_k$ (for Q), $d_k \times N$ (for K^T), d_k (under $\sqrt{d_k}$)

queries: $Q = XW_Q \in \mathbb{R}^{N \times d_k}$

Handwritten note: $d \times d_k$ (under W_Q)

keys: $K = XW_K \in \mathbb{R}^{N \times d_k}$

Handwritten note: $d \times d_k$ (under W_K)

values: $V = XW_V \in \mathbb{R}^{N \times d_v}$

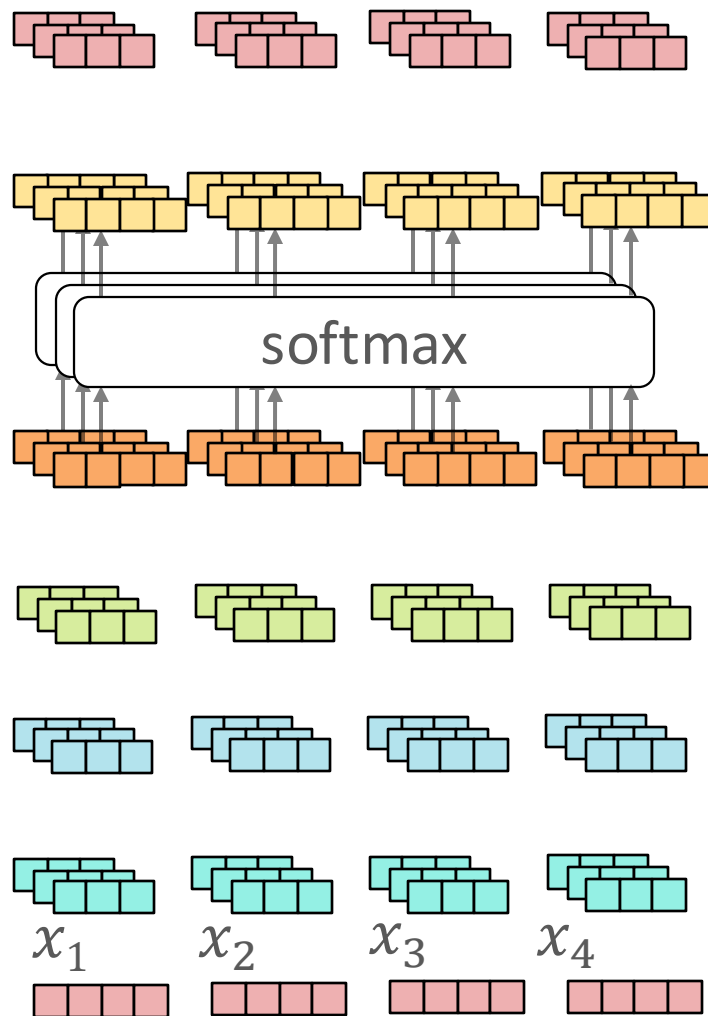
Handwritten note: $d \times d_v$ (under W_V)

design matrix: $X \in \mathbb{R}^{N \times D}$

Handwritten notes: "tokens" (under N), "dim of embedding space" (under D with an arrow)

Multi-head Attention Layer

- Idea: just like we might want multiple convolutional filters in a convolutional layer, we might want multiple attention weights to learn different relationships between tokens!



$$H^{(h)} = \text{softmax}(S^{(h)})V^{(h)}$$

attention weights

$$\text{scores: } S^{(h)} = \frac{Q^{(h)}K^{(h)T}}{\sqrt{d_k^{(h)}}}$$

$$\text{queries: } Q^{(h)} = XW_Q^{(h)}$$

$$\text{keys: } K^{(h)} = XW_K^{(h)}$$

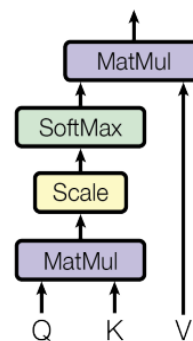
$$\text{values: } V^{(h)} = XW_V^{(h)}$$

design matrix: X

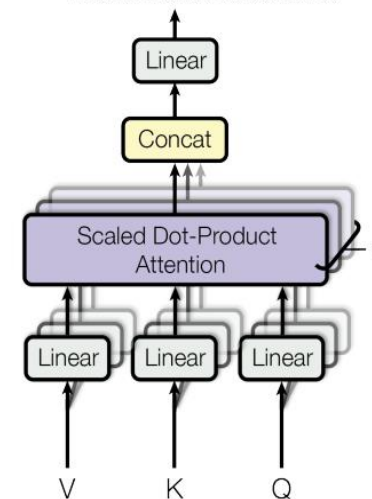
Multi-head Attention Layer

- Idea: just like we might want multiple convolutional filters in a convolutional layer, we might want multiple attention weights to learn different relationships between tokens!

Scaled Dot-Product Attention



Multi-Head Attention



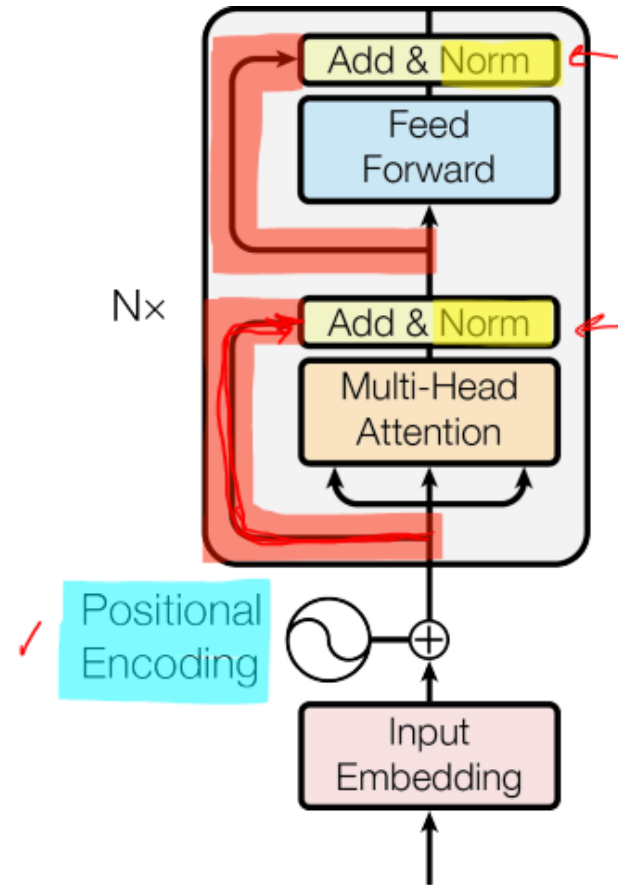
• Suppose h attention heads

- The outputs from all the attention heads are concatenated together to get the final representation

$$H = \left[\underbrace{H^{(1)}}_{d_v}, \underbrace{H^{(2)}}_{d_v}, \dots, \underbrace{H^{(h)}}_{d_v} \right] \quad d_v \times h = D$$

- Common architectural choice: $d_v = D/h \rightarrow |H| = D$

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

Positional Encodings

- Issue: if all tokens attend to every token in the sequence, then how does the model infer the order of tokens?
- Idea: add a position-specific embedding p_t to the token embedding x_t

$$x'_t = x_t + p_t$$

- Positional encodings can be
 - *fixed* i.e., some predetermined function of t or *learned* alongside the token embeddings
 - *absolute* i.e., only dependent on the token's location in the sequence or *relative* to the query token's location

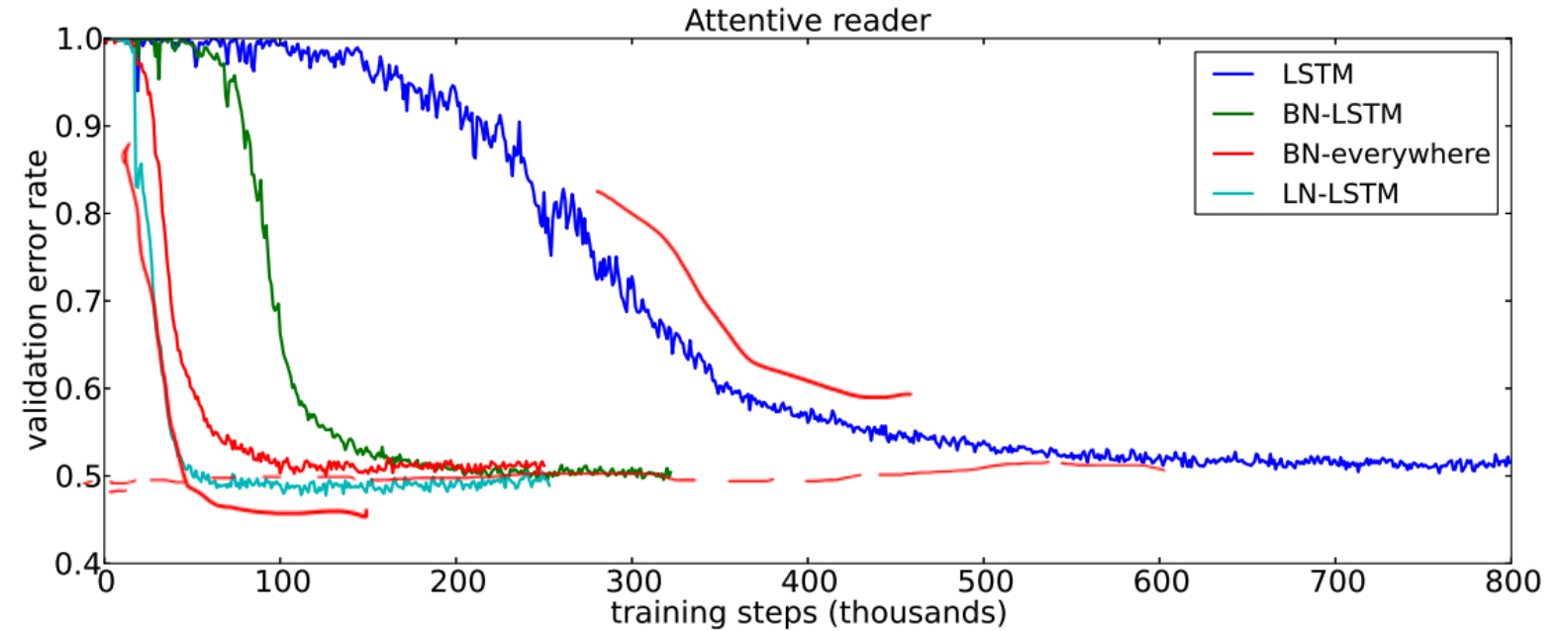
Layer Normalization

- 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, β , and variance, γ^2

$$\underbrace{H'}_{\text{red}} = \gamma \left(\frac{\overbrace{H}^{\text{red}} - \mu}{\sigma} \right) + \beta$$

where μ is the mean and σ 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, β , and variance, γ^2

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

where μ is the mean and σ is the standard deviation of the values in the vector H

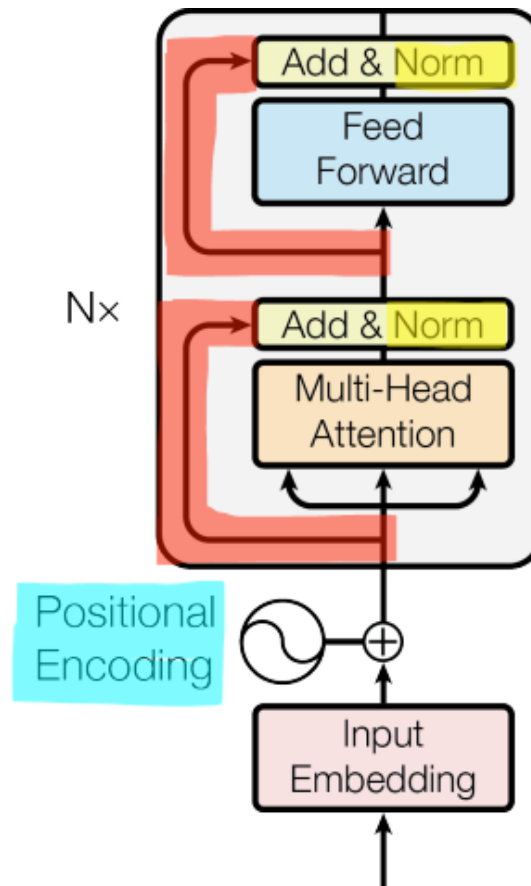
Residual Connections

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

How on earth do we train these things?



- 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

Learning Paradigms

- Supervised learning - $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - Regression - $y^{(i)} \in \mathbb{R}$
 - Classification - $y^{(i)} \in \{1, \dots, C\}$
- Unsupervised learning - $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$
 - Clustering
 - Dimensionality reduction
- Reinforcement learning
- Active learning
- Semi-supervised learning
- Online learning

Learning Paradigms

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

Clustering

$$D = \{x^{(i)}\}_{i=1}^N$$

- 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:

Recall: Similarity for k NN

- Classify a point as the label of the “most similar” training point
- **Idea: given real-valued features, we can use a distance metric to determine how similar two data points are**
- A common choice is Euclidean distance:

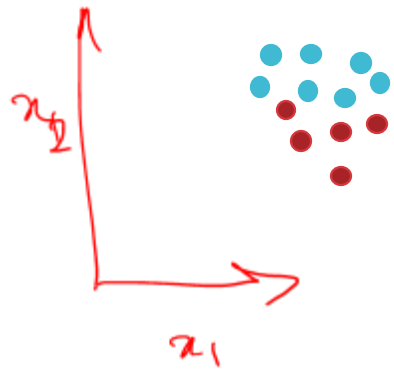
$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_2 = \sqrt{\sum_{d=1}^D (x_d - x'_d)^2}$$

- An alternative is the Manhattan distance:

$$d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_1 = \sum_{d=1}^D |x_d - x'_d|$$

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 partition is best?

General Recipe for Machine Learning

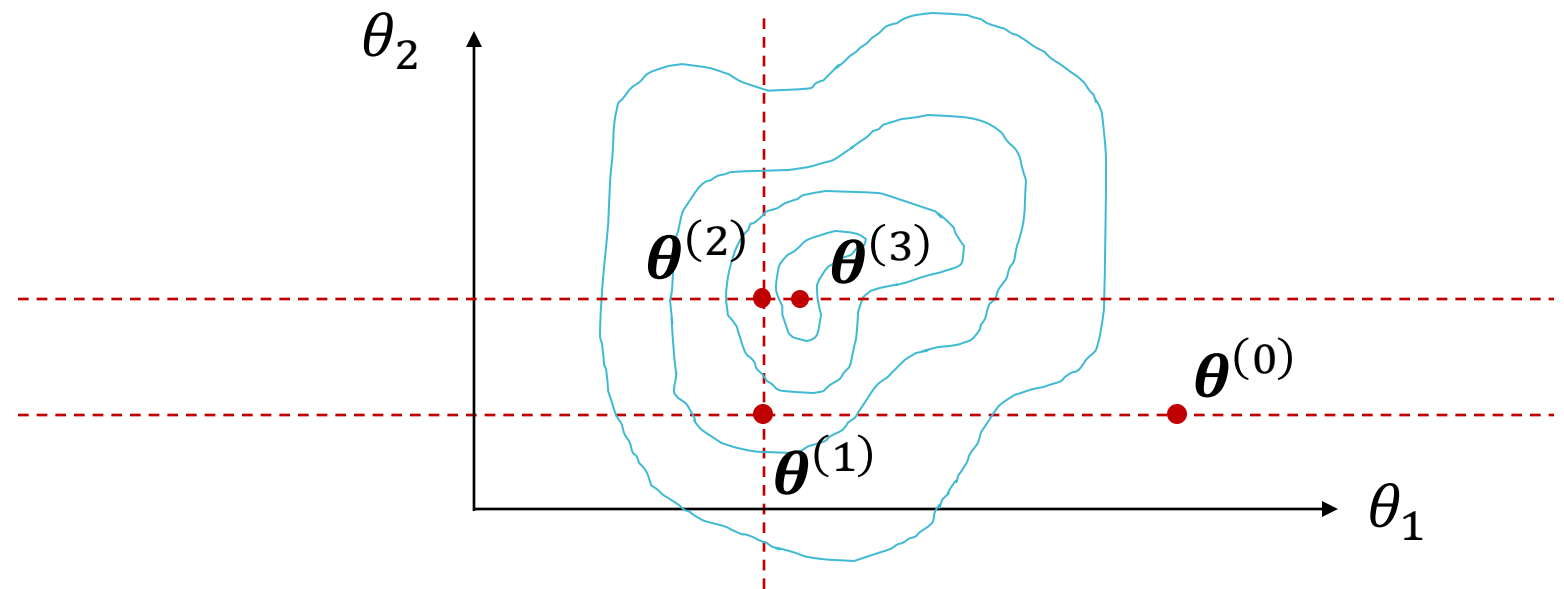
- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

Recipe for K -means

- Define a model and model parameters $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$
 - Assume K clusters and use the Euclidean distance
 - Parameters: μ_1, \dots, μ_K and $z^{(1)}, \dots, z^{(N)}$
center of cluster 1, ..., K *cluster assigned to $x^{(i)}$*
- Write down an objective function
$$\rightarrow \sum_{i=1}^N \|x^{(i)} - \mu_{z^{(i)}}\|_2$$
- Optimize the objective w.r.t. the model parameters
 - Use (block) coordinate descent

Coordinate Descent

- Goal: minimize some objective $\vec{\theta} = (\theta_1, \theta_2)$
 $\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta)$
- 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

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

- If μ_1, \dots, μ_K are fixed

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

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

$$\begin{aligned} \hat{\mu}_k &= \operatorname{argmin}_{\mu} \sum_{i: z^{(i)} = k} \|\mathbf{x}^{(i)} - \mu\|_2 \\ &= \frac{1}{N_k} \sum_{i: z^{(i)} = k} \mathbf{x}^{(i)} \end{aligned}$$

K-means Algorithm

- Input: $\mathcal{D} = \{(\mathbf{x}^{(i)})\}_{i=1}^N, K$
- 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:

$$z^{(i)} = \underset{k}{\operatorname{argmin}} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_k\|_2$$

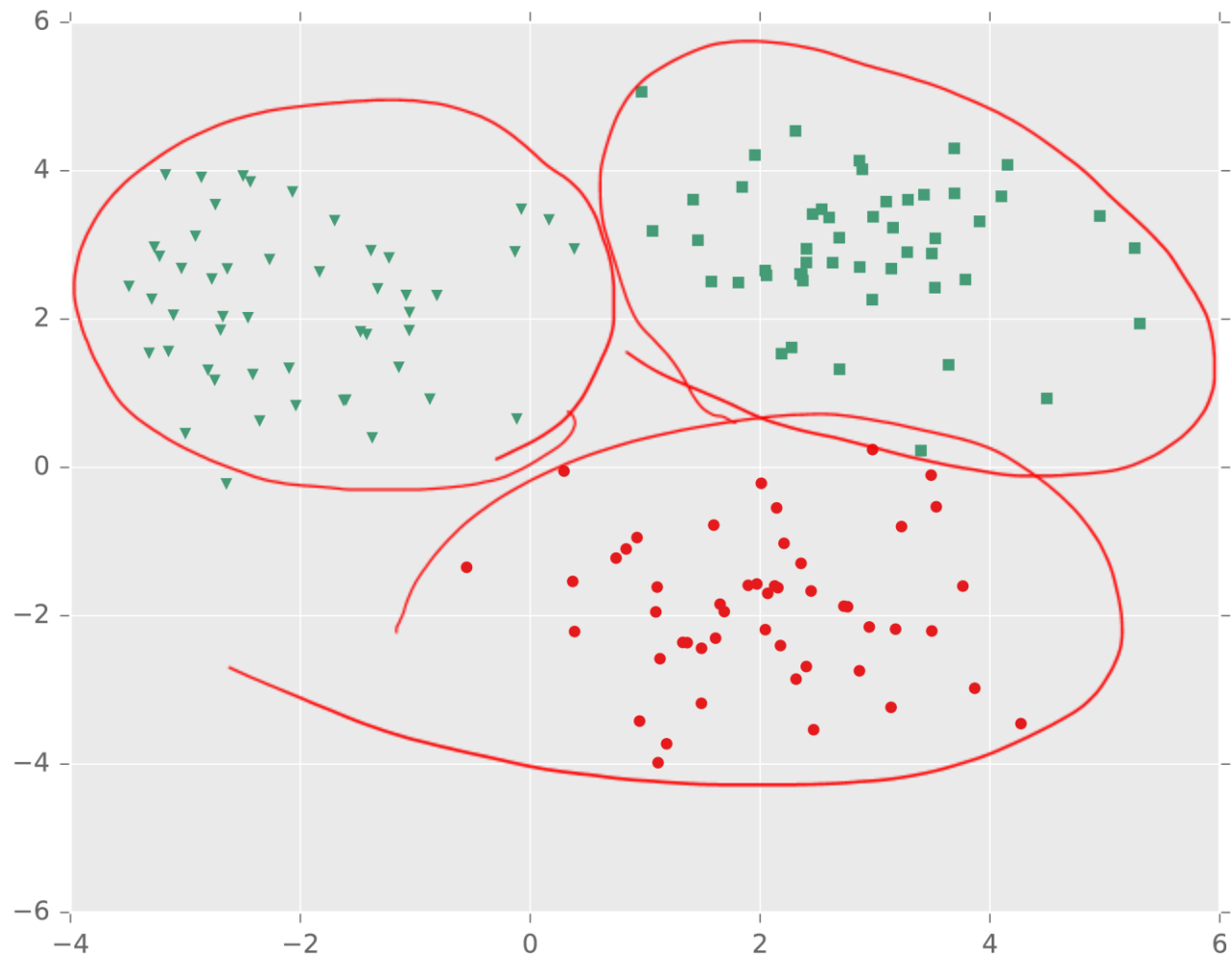
- b. Recompute the cluster centers:

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i: z^{(i)}=k} \mathbf{x}^{(i)}$$

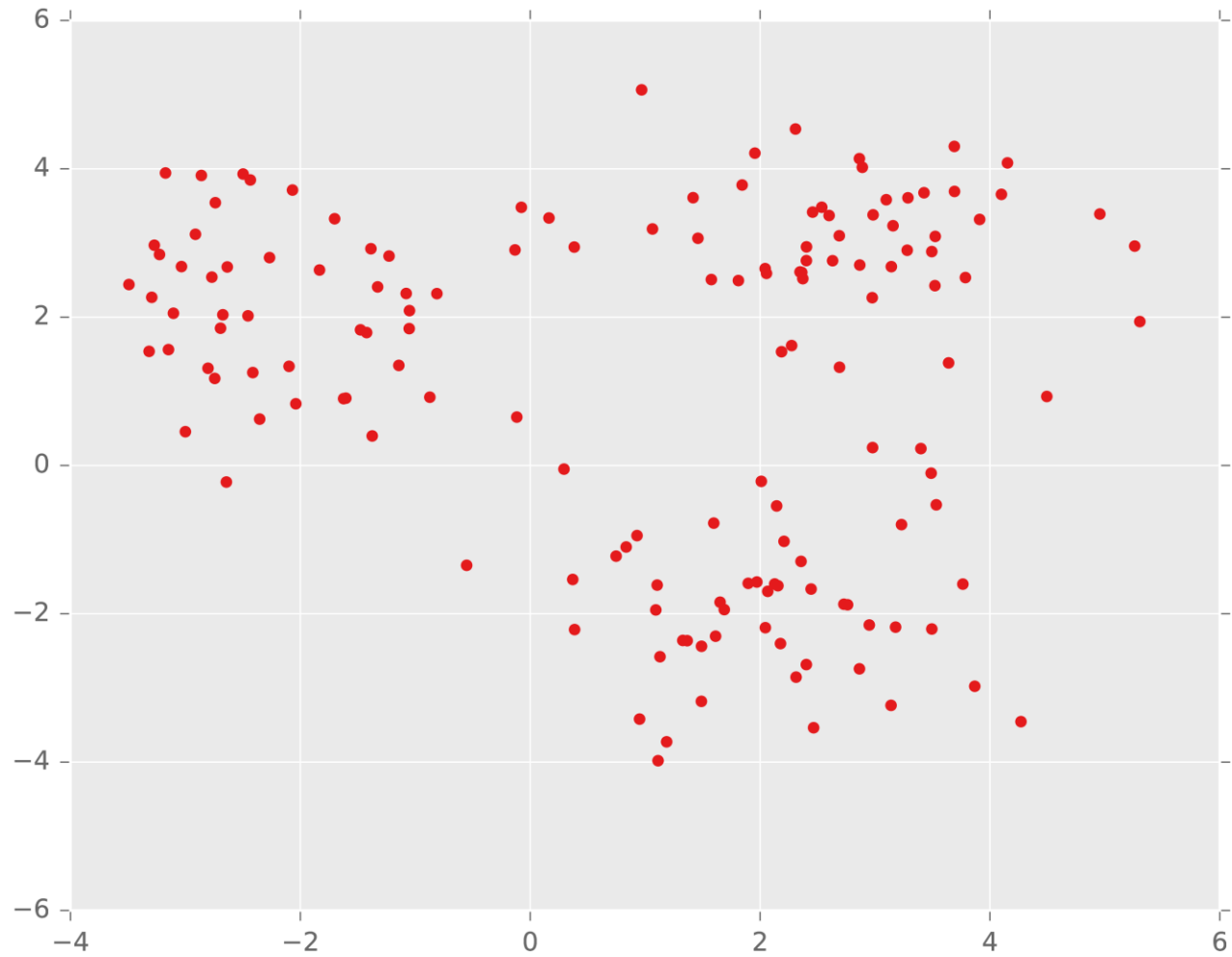
where N_k is the number of data points in cluster k

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

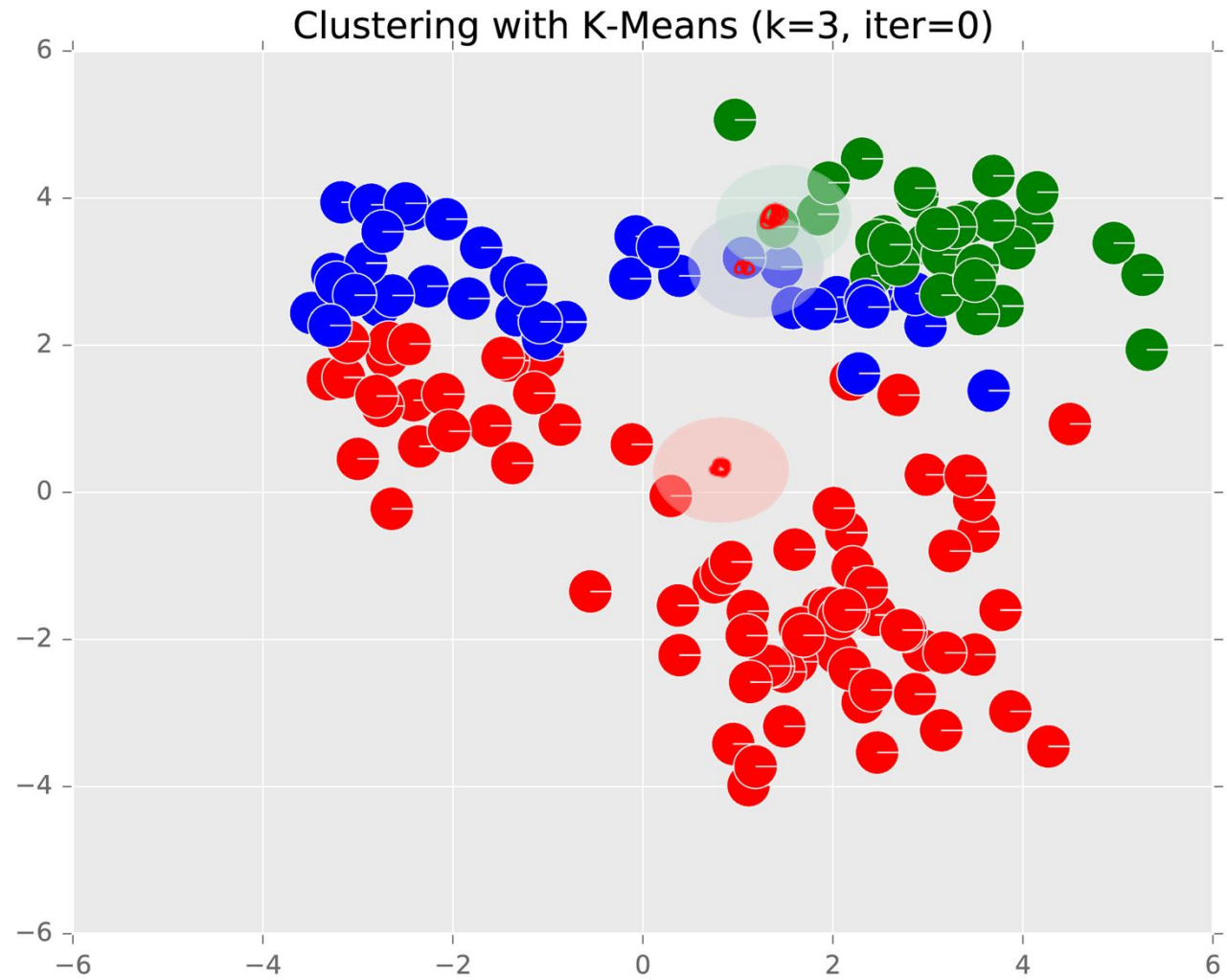
K -means: Example ($K = 3$)



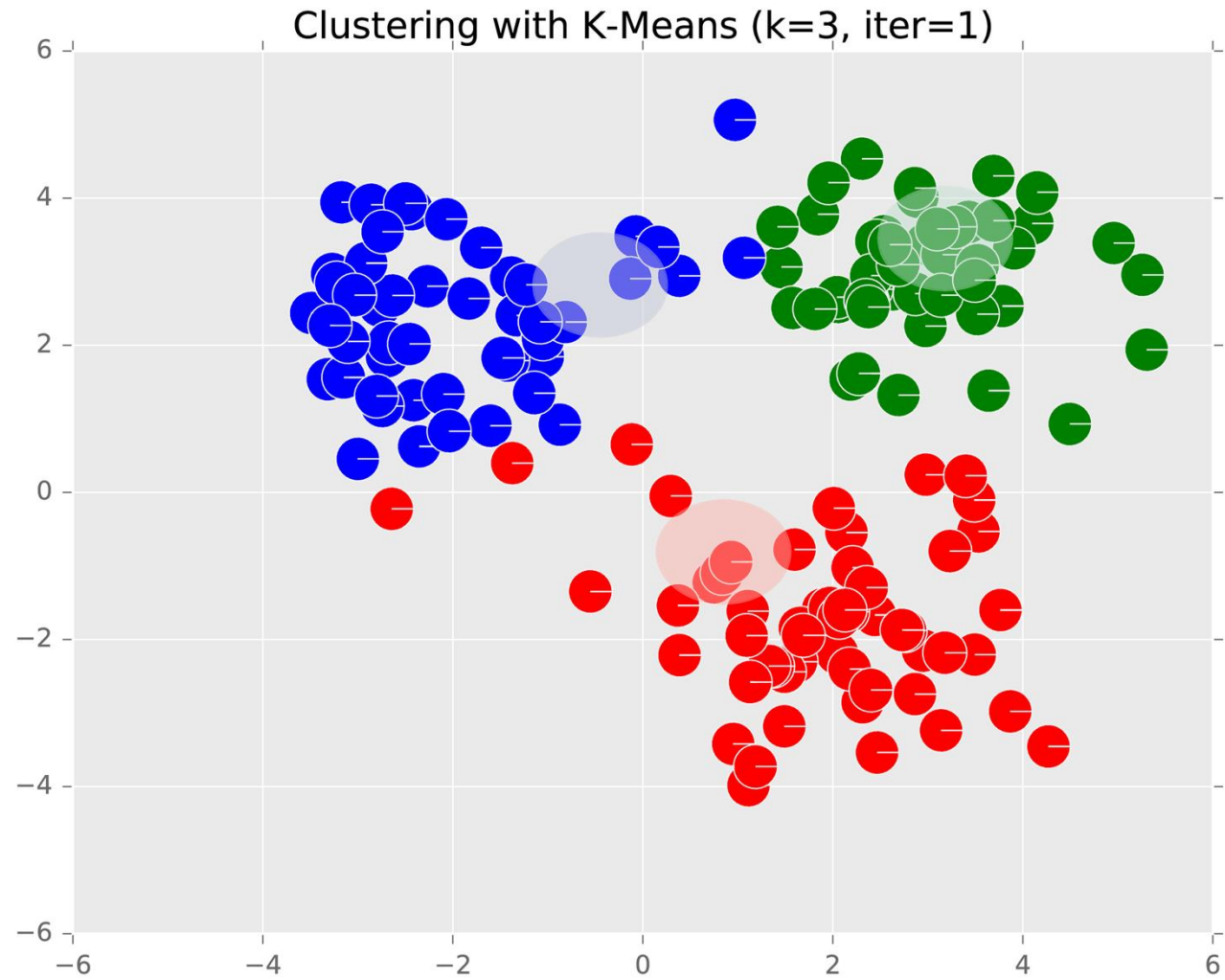
K -means: Example ($K = 3$)



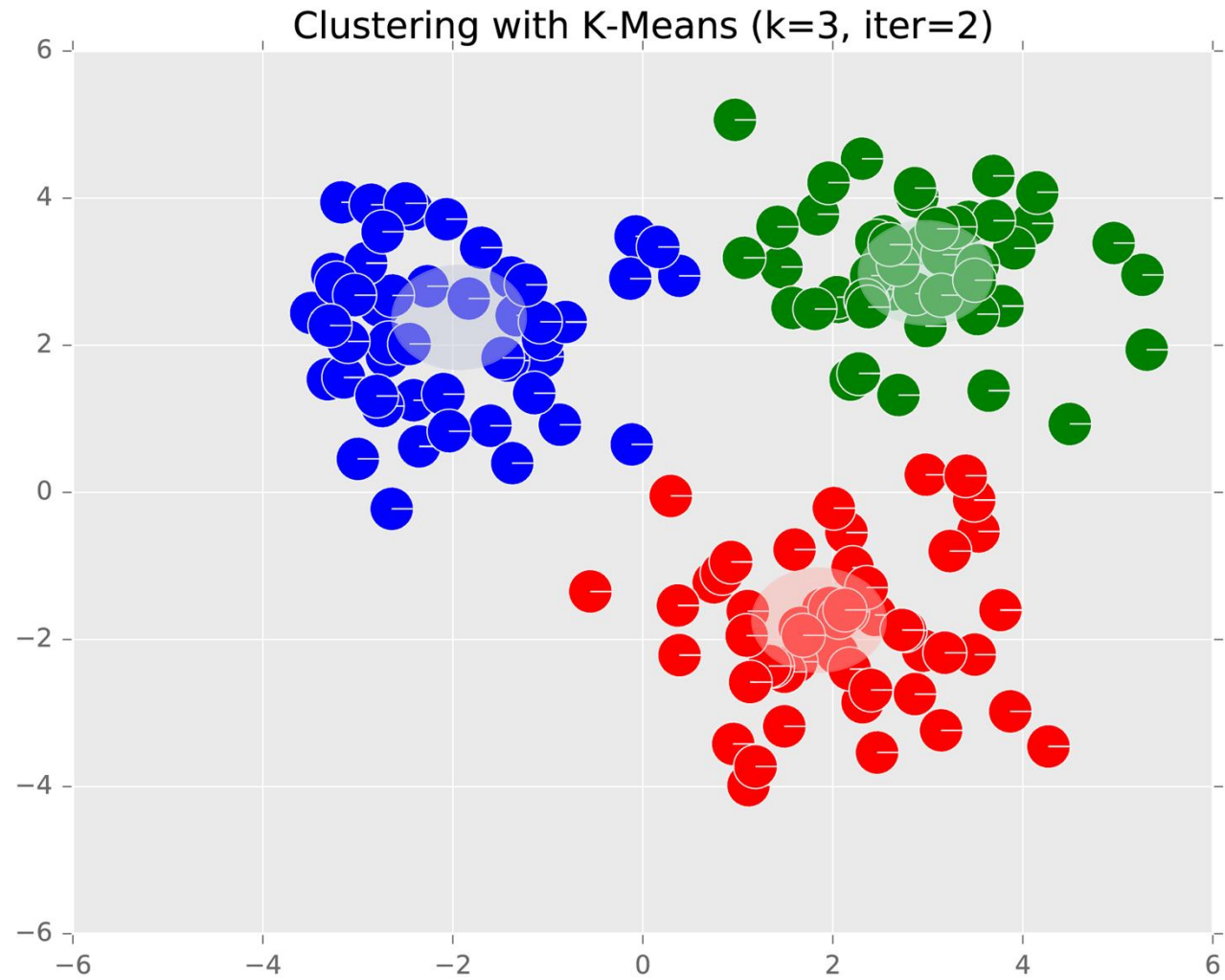
K -means: Example ($K = 3$)



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



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



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



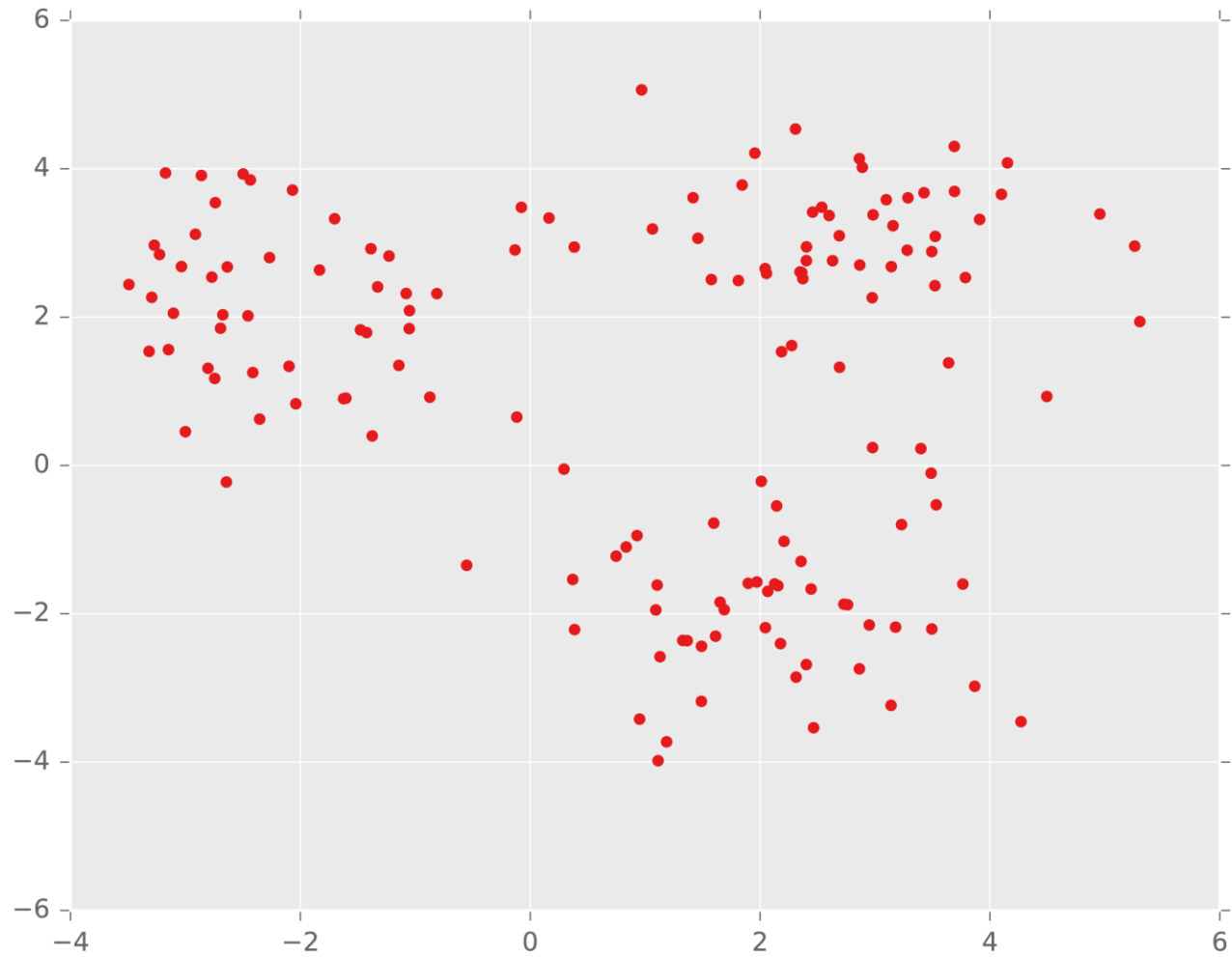
K -means: Example ($K = 3$)



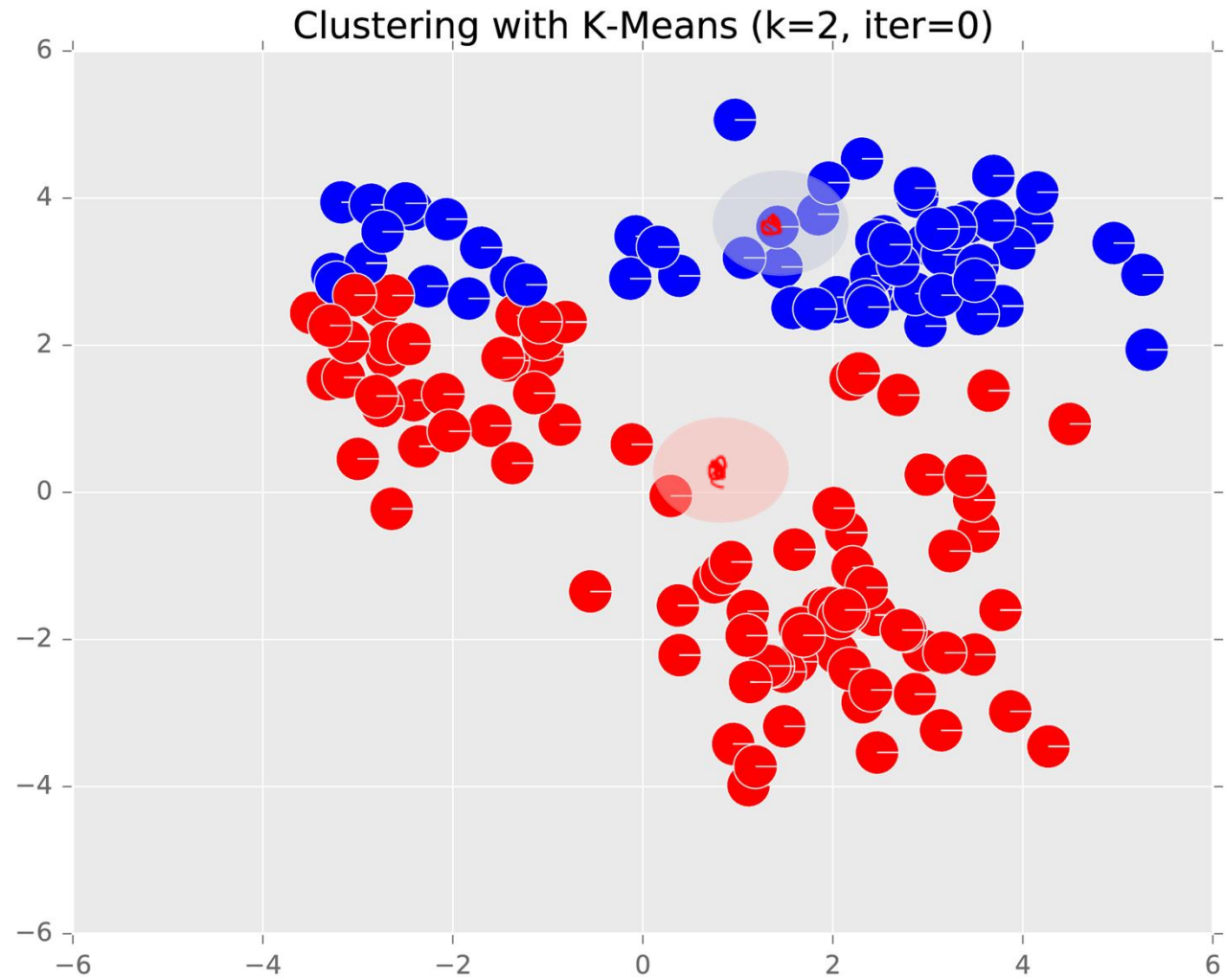
K -means: Example ($K = 3$)



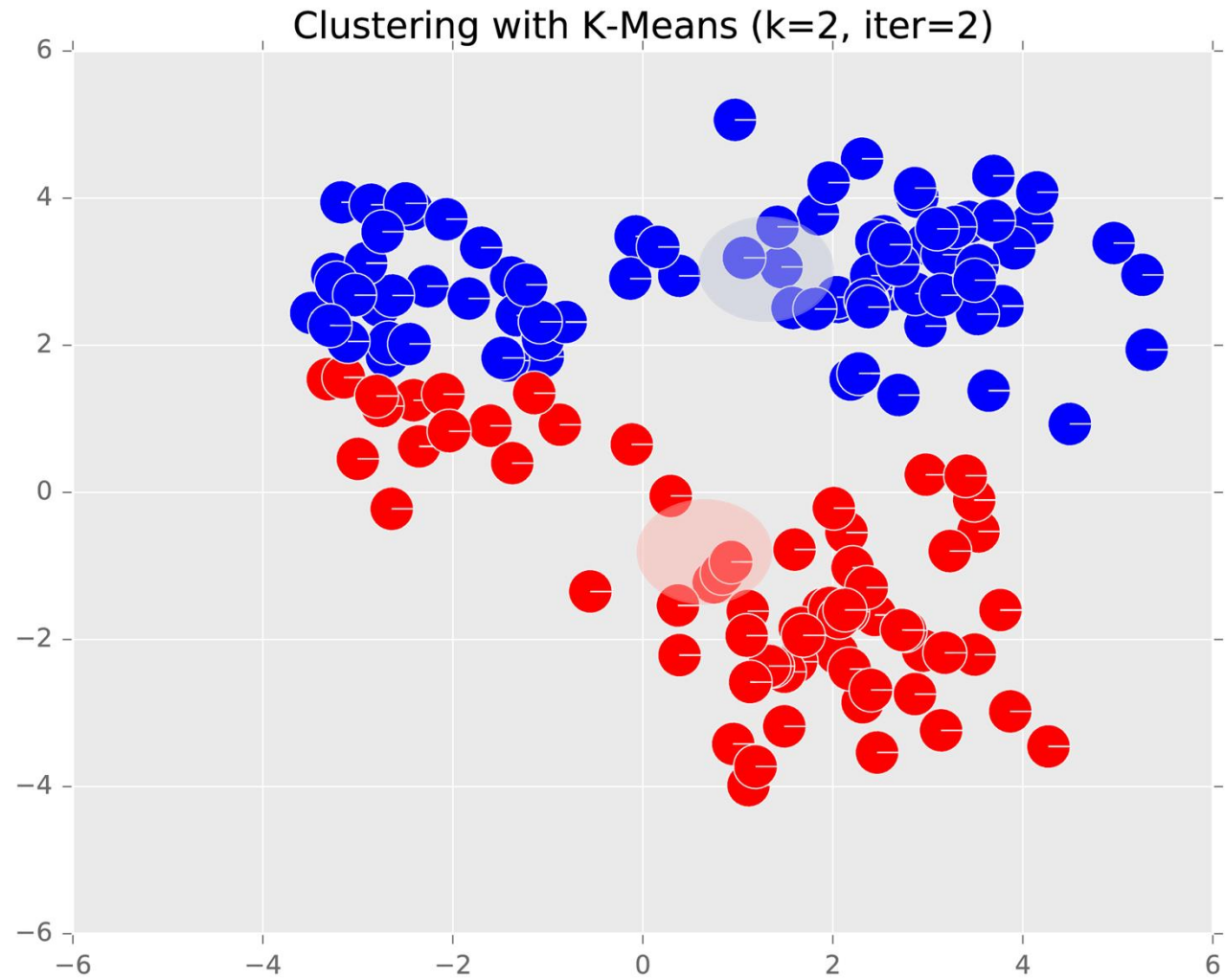
K -means: Example ($K = 2$)



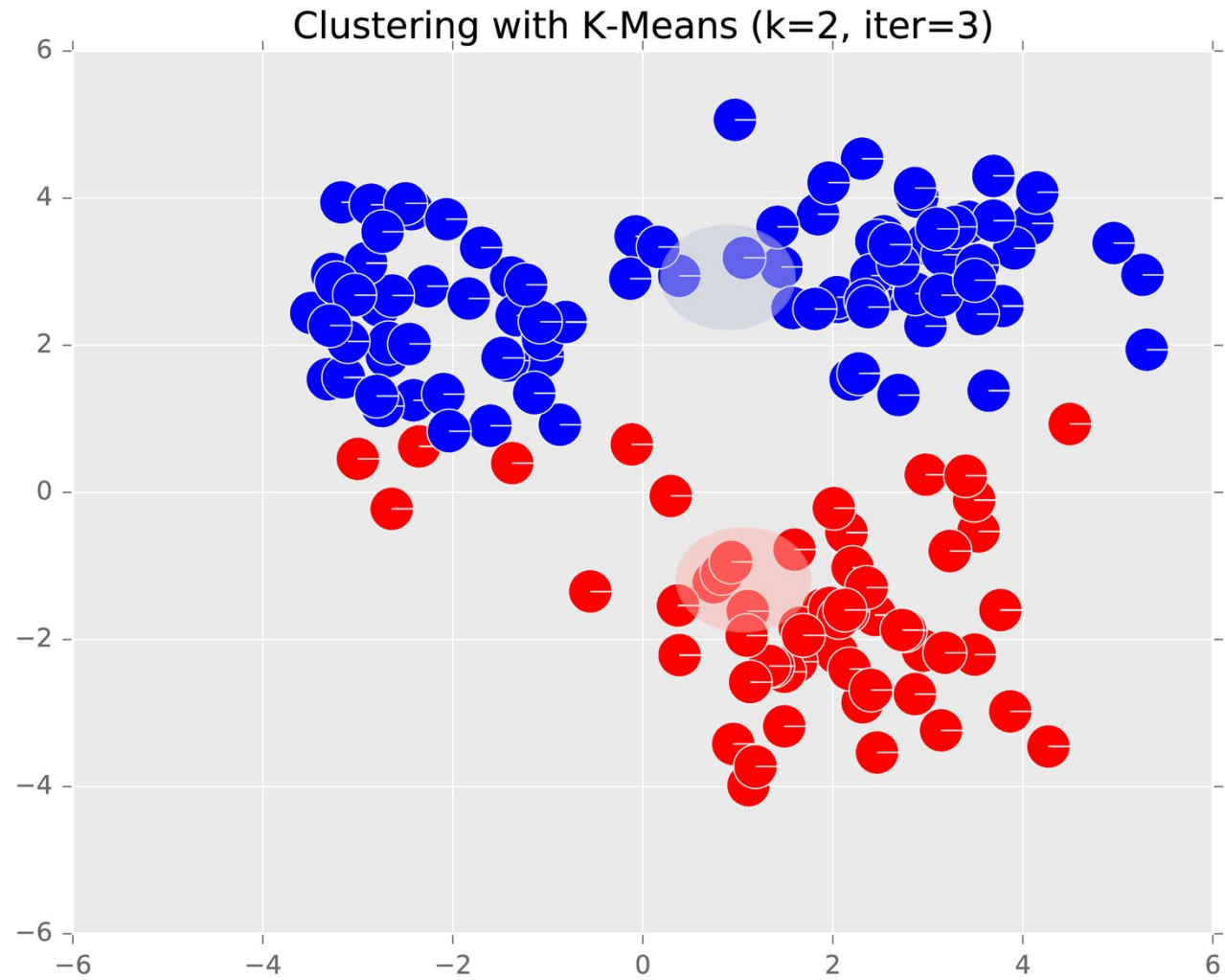
K -means: Example ($K = 2$)



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



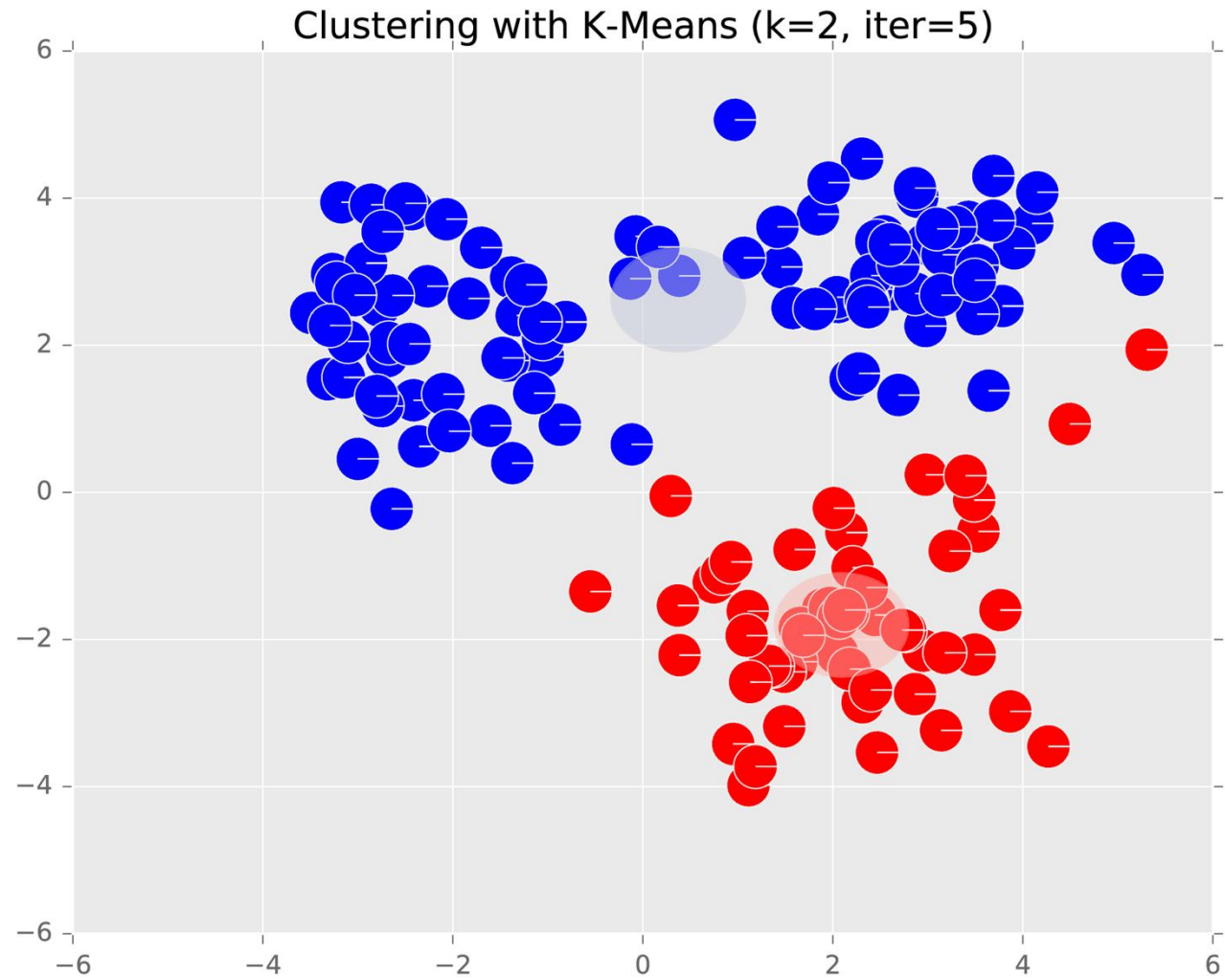
K -means: Example ($K = 2$)



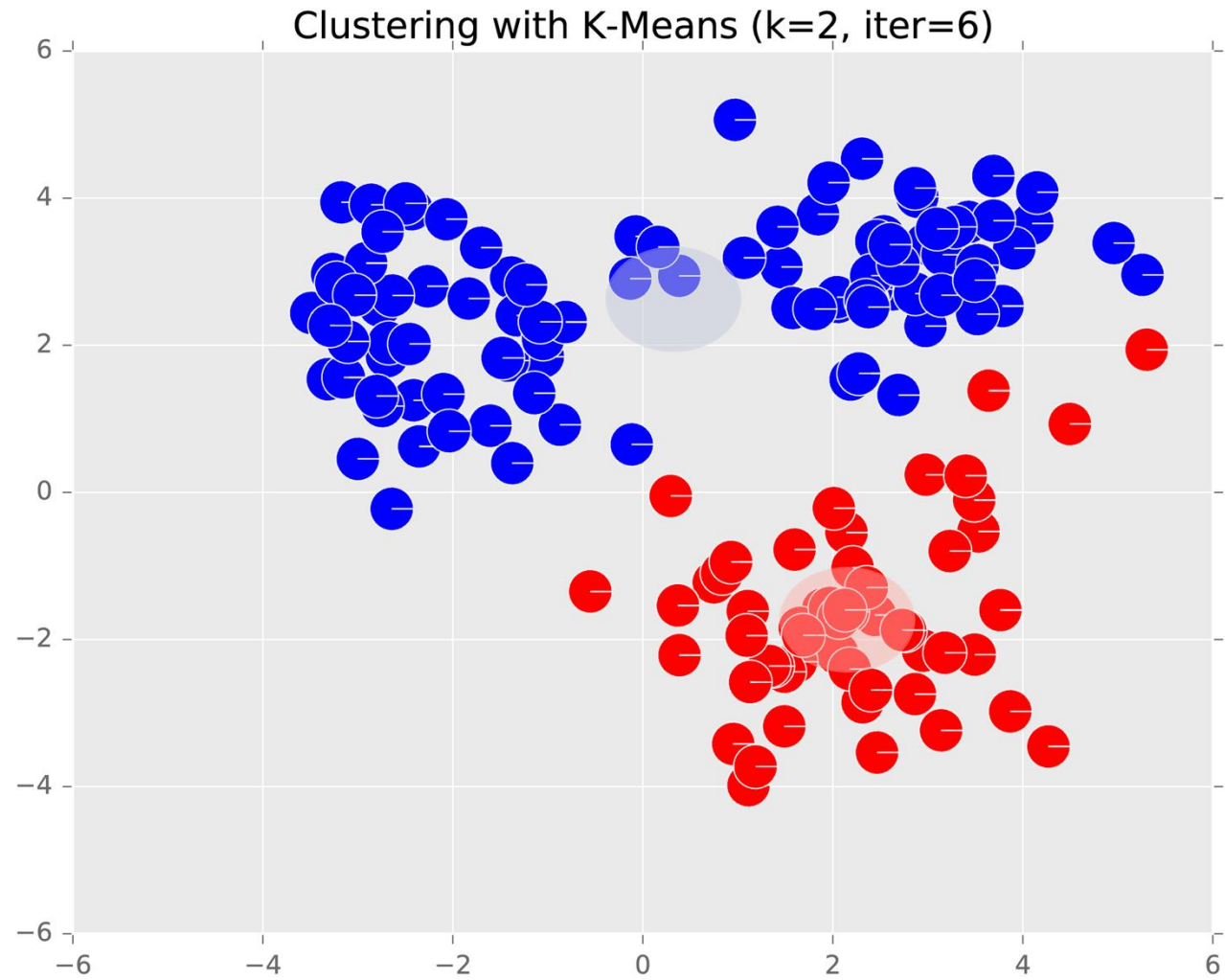
K -means: Example ($K = 2$)



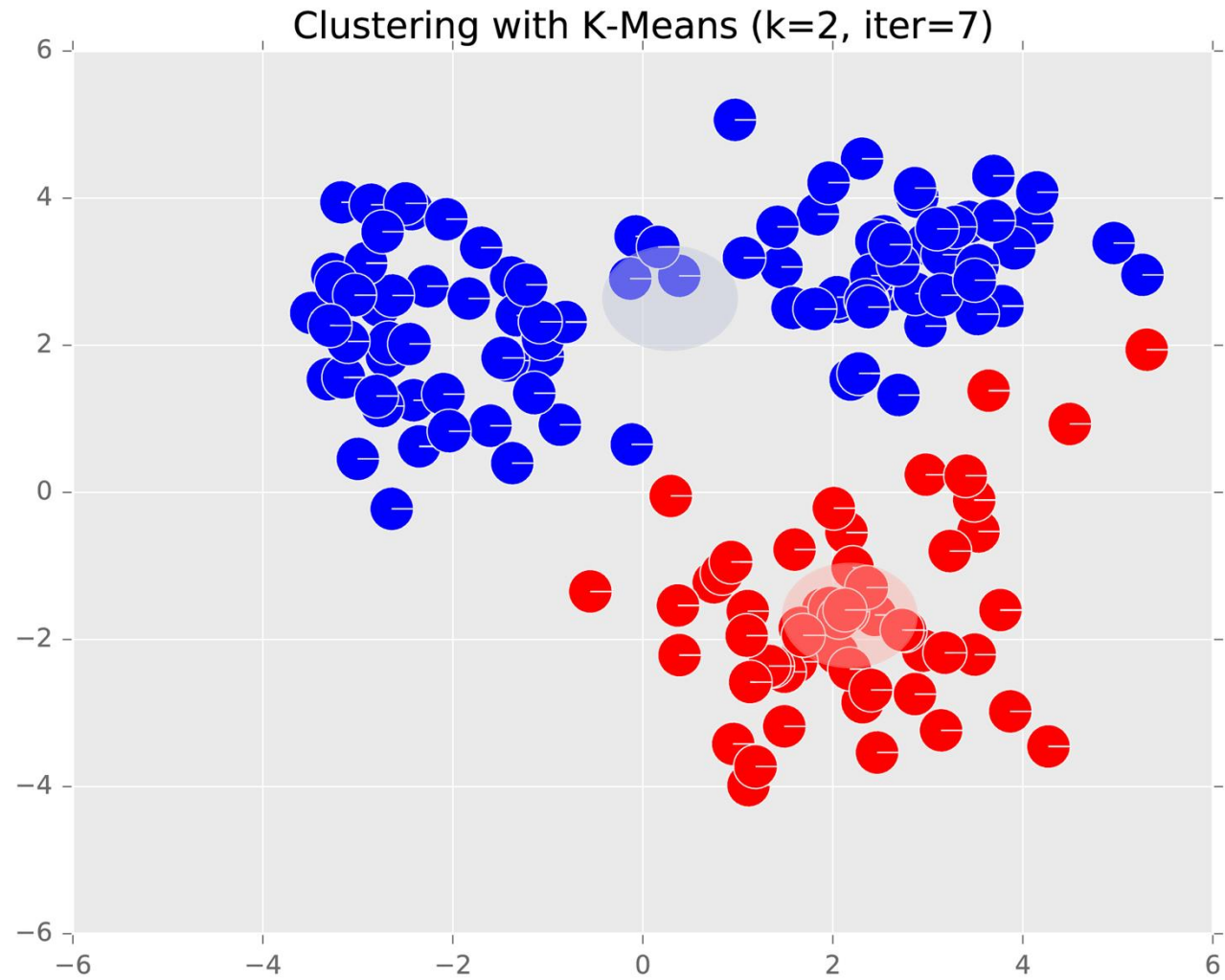
K -means: Example ($K = 2$)



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



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

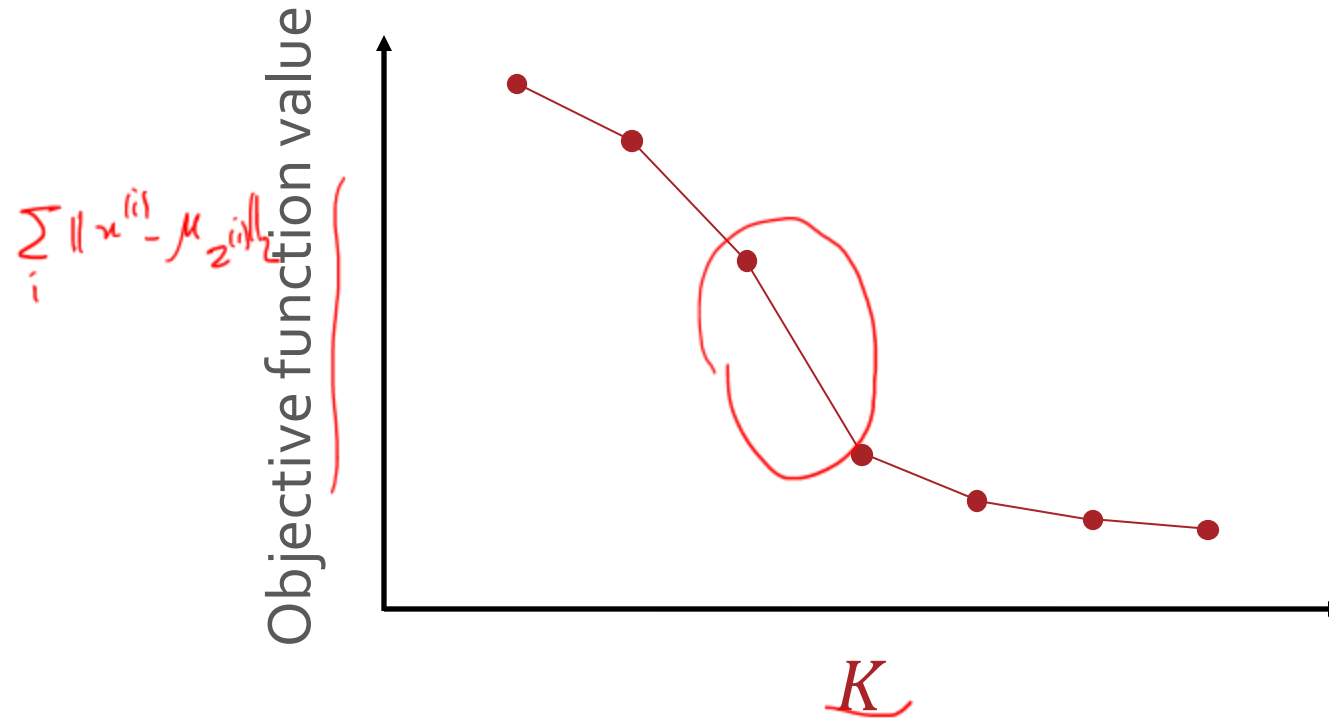


Setting K

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

Setting K

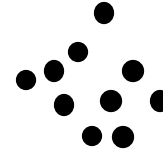
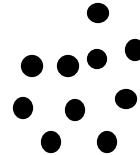
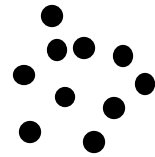
- Idea: choose the value of K that minimizes the objective function



- Better idea: look for the characteristic “elbow” or largest decrease when going from $K - 1$ to K

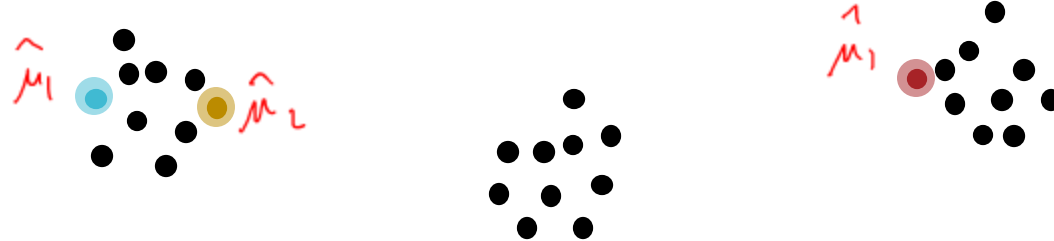
Initializing K -means

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



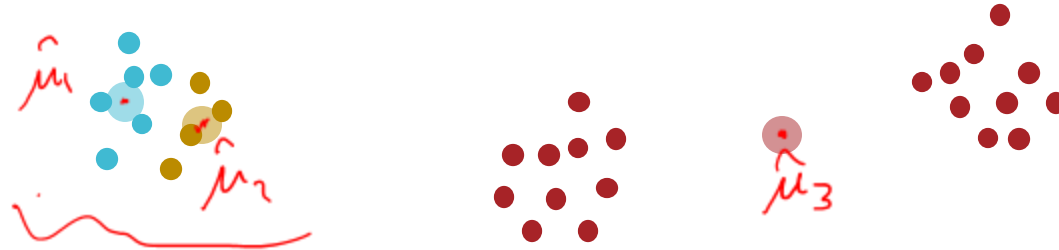
Initializing K -means

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



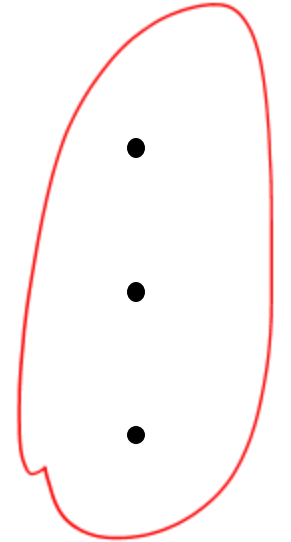
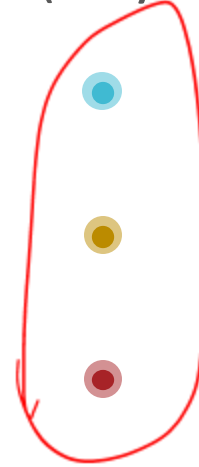
Initializing K -means

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



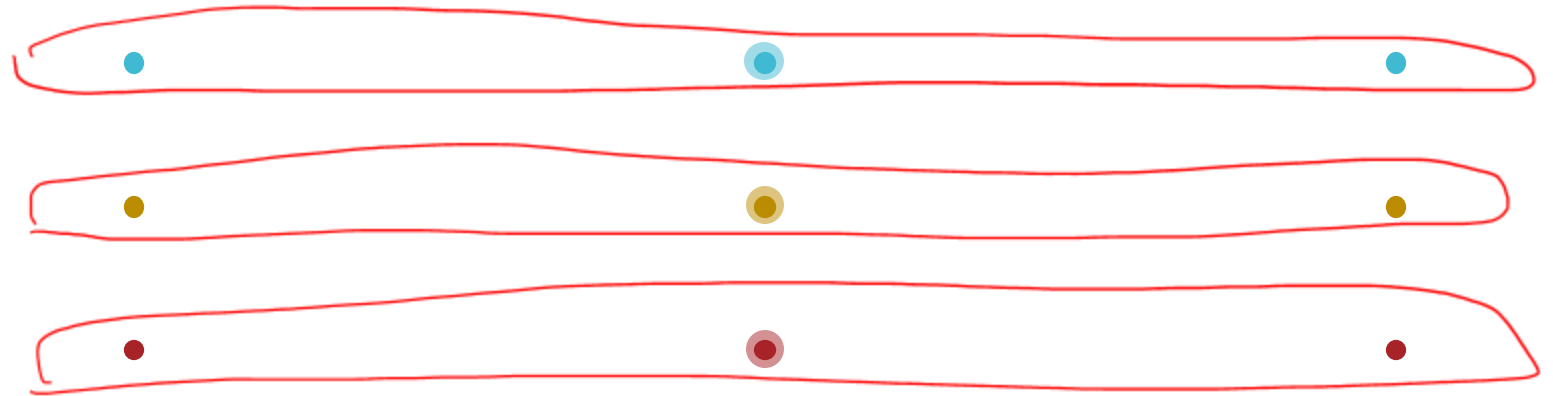
Initializing K -means

- 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)

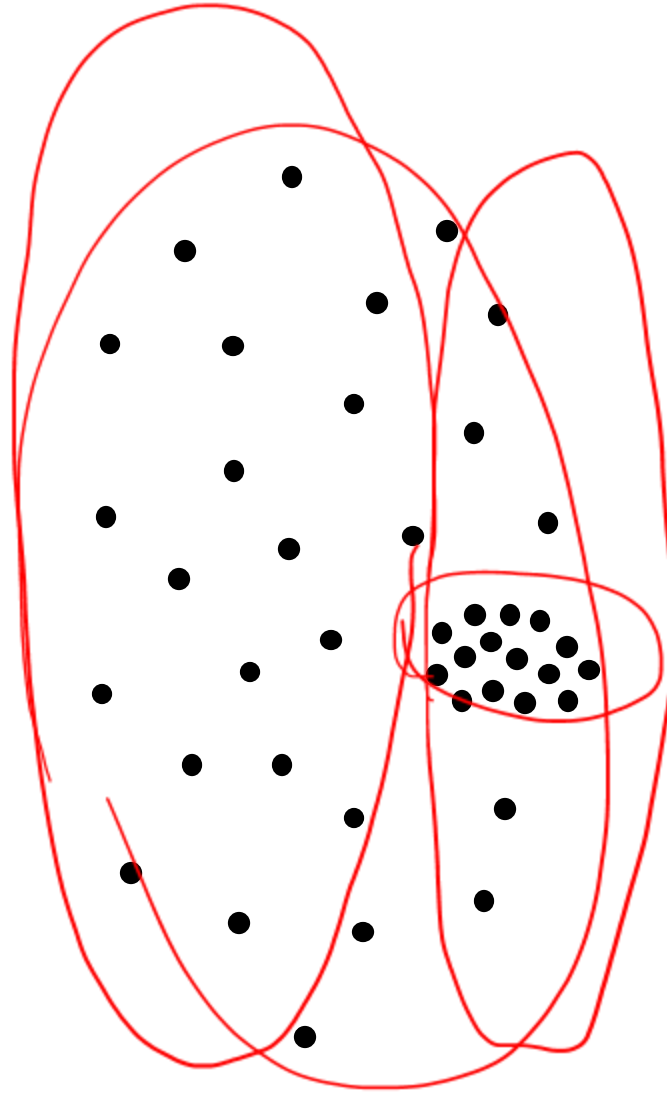


- Lloyd's method converges to a local minimum and that local minimum can be arbitrarily bad (relative to the optimal clusters)
- Intuition: want initial cluster centers to be far apart from one another

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.

Shortcomings of K -means



- Clusters cannot overlap
- Clusters must all be of the same “width”
- Clusters must be linearly separable

Probabilistic or “Soft” Assignments

- Instead of $z^{(i)}$ being a deterministic scalar, let $\mathbf{z}^{(i)}$ be a 1-of- K vector indicating cluster membership
 - For example, $\mathbf{z}^{(1)} = [0, 1, 0, \dots, 0]$ indicates that the first data point belongs to the second cluster
- Let $\pi_k := p\left(z_k^{(i)} = 1\right)$

Gaussian Mixture Models (GMMs)

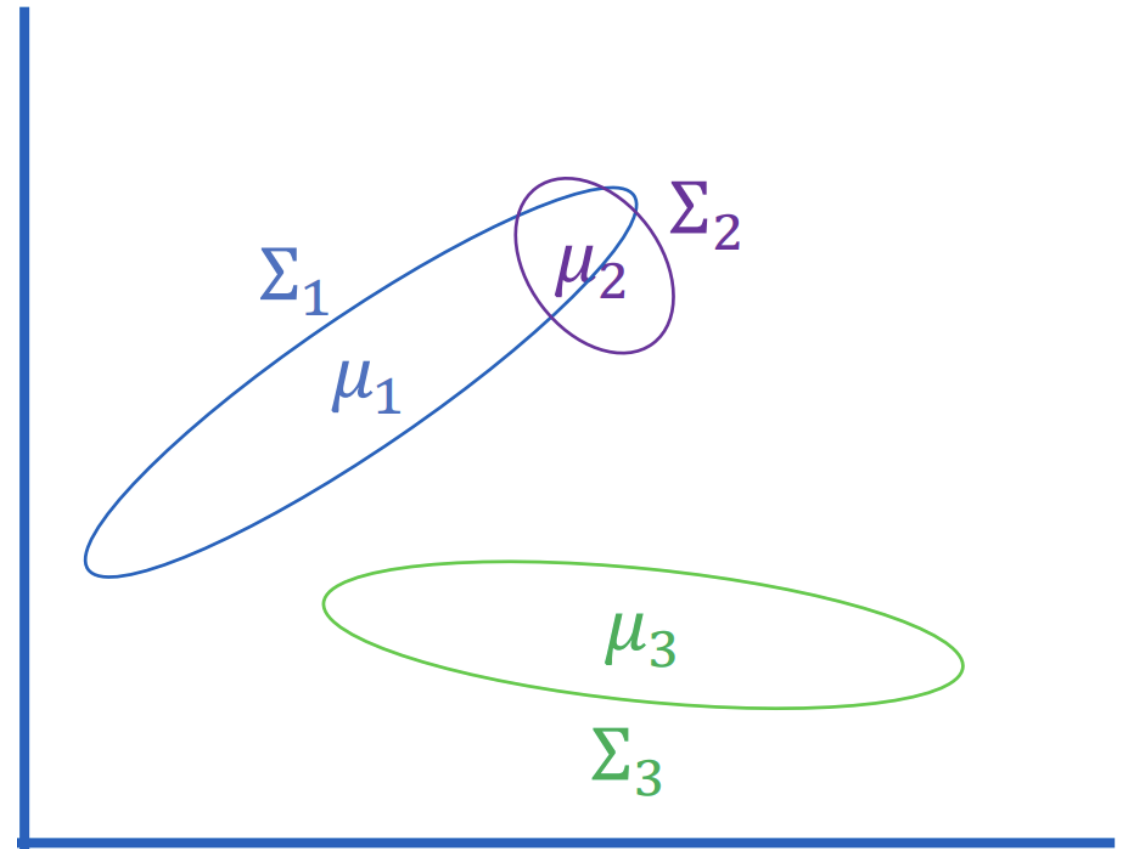
Assume the following data-generating model for our dataset, $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$

1. Sample a cluster at random:

$$p(z_k^{(i)} = 1) = \pi_k$$

2. Sample a data point from the chosen cluster:

$$p(\mathbf{x}^{(i)} | z_k^{(i)} = 1) \sim N(\mu_k, \Sigma_k)$$



Gaussian Mixture Models (GMMs)

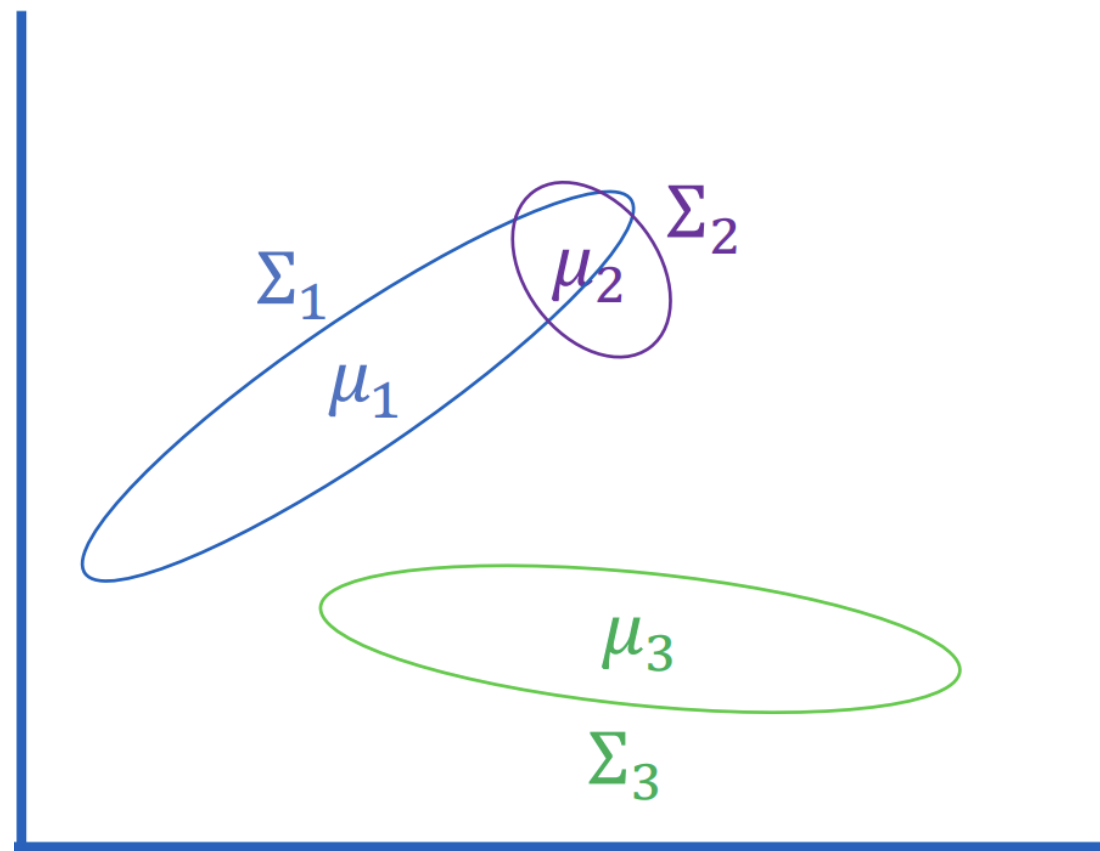
Assume the following data-generating model for our dataset, $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$

1. Sample a cluster at random:

$$p(z_k^{(i)} = 1) = \pi_k$$

2. Sample a data point from the chosen cluster:

$$p(\mathbf{x}^{(i)} | z_k^{(i)} = 1) \sim N(\mu_k, \Sigma_k)$$



Let $\theta = \{\mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K, \pi_1, \dots, \pi_K\}$

Maximizing the Likelihood?

- The log likelihood of $\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}\}_{i=1}^N$ is $\ell_{\mathcal{D}}(\theta) =$

Maximizing the Likelihood?

- The log likelihood of $\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}\}_{i=1}^N$ is

$$\begin{aligned}\ell_{\mathcal{D}}(\theta) &= \log \prod_{i=1}^N p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) \\&= \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}, \theta) + \log p(\mathbf{z}^{(i)} | \theta) \\&= \sum_{i=1}^N \log \prod_{k=1}^K p(\mathbf{x}^{(i)} | z_k^{(i)} = 1, \theta)^{z_k^{(i)}} + \log \prod_{k=1}^K p(z_k^{(i)} = 1 | \theta)^{z_k^{(i)}} \\&= \sum_{i=1}^N \sum_{k=1}^K z_k^{(i)} \log p(\mathbf{x}^{(i)} | z_k^{(i)} = 1, \theta) + \sum_{k=1}^K z_k^{(i)} \log p(z_k^{(i)} = 1 | \theta) \\&= \sum_{i=1}^N \sum_{k=1}^K z_k^{(i)} (\log N(\mathbf{x}^{(i)}; \mu_k, \Sigma_k) + \log \pi_k)\end{aligned}$$

Maximizing
the
Complete
Likelihood is
easy but
requires $z^{(i)}$!

- The log complete likelihood of $\mathcal{D} = \{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}\}_{i=1}^N$ is

$$\begin{aligned}\ell_{\mathcal{D}}(\theta) &= \log \prod_{i=1}^N p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \theta) \\&= \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}, \theta) + \log p(\mathbf{z}^{(i)} | \theta) \\&= \sum_{i=1}^N \log \prod_{k=1}^K p(\mathbf{x}^{(i)} | z_k^{(i)} = 1, \theta)^{z_k^{(i)}} + \log \prod_{k=1}^K p(z_k^{(i)} = 1 | \theta)^{z_k^{(i)}} \\&= \sum_{i=1}^N \sum_{k=1}^K z_k^{(i)} \log p(\mathbf{x}^{(i)} | z_k^{(i)} = 1, \theta) + \sum_{k=1}^K z_k^{(i)} \log p(z_k^{(i)} = 1 | \theta) \\&= \sum_{i=1}^N \sum_{k=1}^K z_k^{(i)} (\log N(\mathbf{x}^{(i)}; \mu_k, \Sigma_k) + \log \pi_k)\end{aligned}$$

- Parameters decoupled \rightarrow set partial derivatives equal to 0

Maximizing the Marginal Likelihood

- The log *marginal* likelihood of $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ is

$$\begin{aligned}\ell(\theta|\mathcal{D}) &= \log \prod_{i=1}^N p(\mathbf{x}^{(i)}|\theta) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}|\theta) \\ &= \sum_{i=1}^N \log \sum_{\mathbf{z}^{(i)}} p(\mathbf{x}^{(i)}|\mathbf{z}^{(i)}, \theta) p(\mathbf{z}^{(i)}|\theta) \\ &= \sum_{i=1}^N \log \sum_{\mathbf{z}^{(i)}} \prod_{k=1}^K \left(p(\mathbf{x}^{(i)}|z_k^{(i)} = 1, \theta) p(z_k^{(i)} = 1|\theta) \right)^{z_k^{(i)}} \\ &= \sum_{i=1}^N \log \sum_{\mathbf{z}^{(i)}} \prod_{k=1}^K \left(N(\mathbf{x}^{(i)}; \mu_k, \Sigma_k) \pi_k \right)^{z_k^{(i)}}\end{aligned}$$

- Parameters coupled and *constrained* \rightarrow gradient ascent is possible but complicated and slow to converge

Recipe for GMMs

- Define a model and model parameters
 - Assume K Gaussian clusters
 - Parameters: $\theta = \{\mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K, \pi_1, \dots, \pi_K\}$
 - Write down an objective function
 - Maximize the log marginal likelihood
- $$\ell_{\mathcal{D}}(\theta) = \log \prod_{i=1}^N p(\mathbf{x}^{(i)} | \theta)$$
- Optimize the objective w.r.t. the model parameters
 - Expectation-maximization

Expectation- Maximization for GMMs: Intuition

- Insight: if we knew the cluster assignments, $\mathbf{z}^{(i)}$, we could maximize the log complete likelihood instead of the log marginal likelihood
- Idea: replace $\mathbf{z}^{(i)}$ in the log complete likelihood with our “best guess” for $\mathbf{z}^{(i)}$ given the parameters and the data
- Observation: changing the parameters changes our “best guess” and vice versa
- Approach: iterate between updating our “best guess” and updating the parameters