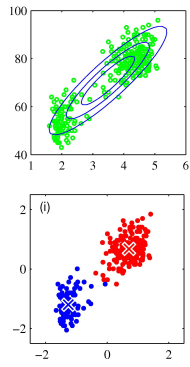


## K-Means Clustering

6

## The K-Means Algorithm

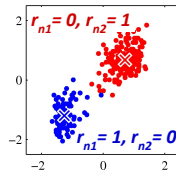
- Given **unlabeled** data  $\{\mathbf{x}^{(n)}\} (n = 1, \dots, N)$ ,
- and believing it belongs in  $K$  clusters (say  $K = 2$  here),
- How do we find the clusters?
  - What would be the objective function?



7

## The K-Means Algorithm

- Use indicator variables  $r_{nk} \in \{0, 1\}$ :
  - $r_{nk} = 1$  if  $\mathbf{x}^{(n)}$  is in cluster  $k$
  - and  $r_{nk} = 0$  for all  $j \neq k$
- Find cluster centers  $\mu_k$  and assignments  $r_{nk}$  to minimize the distortion measure  $J$ :
  - Sum of squared distance of points from the center of its own cluster (*Intra-cluster variation*):



$$J = \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|\mathbf{x}^{(n)} - \mu_k\|^2 \quad \mu_k = \frac{1}{N_k} \sum_{n: \mathbf{x}^{(n)} \in \text{cluster } k} \mathbf{x}^{(n)} = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}^{(n)}}{\sum_{n=1}^N r_{nk}}$$

8

## The K-Means Algorithm

- Initialize the cluster centers (centroids)
- Repeat the following update until convergence:
  - $r := \arg \min_r J(r, \mu)$
  - $\mu := \arg \min_\mu J(r, \mu)$

$$\text{where } J = \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|\mathbf{x}^{(n)} - \mu_k\|^2$$

$$\mu_k = \frac{1}{N_k} \sum_{n: \mathbf{x}^{(n)} \in \text{cluster } k} \mathbf{x}^{(n)} = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}^{(n)}}{\sum_{n=1}^N r_{nk}}$$

9

## The K-Means Algorithm

- Initialize the cluster centers.
- Repeat until convergence:
  - Cluster assignment ("E-Step"):** assign each point to closest center.
 
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}^{(n)} - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$
  - Parameter update ("M-Step"):** update the centers

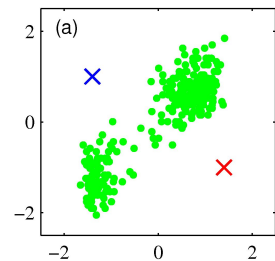
Note: E, M stands for:  
 • E: Expectation  
 • M: Maximization  
 (We will revisit EM later.)

$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}^{(n)}}{\sum_n r_{nk}} \quad \text{Q. Verify this}$$

10

## K-Means Clustering

- Select K. Pick random centroids.
  - Here  $K=2$ .

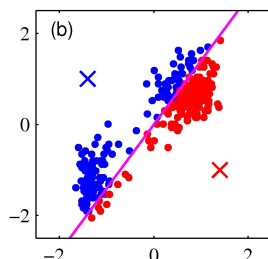


11

## K-Means Clustering

### Cluster assignment Step ("E-Step")

- Assign each point to the nearest center.

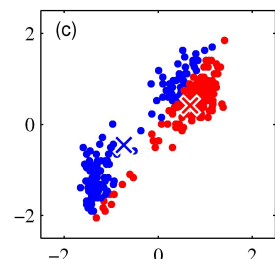


12

## K-Means Clustering

### Update parameters (centroids) ("M-Step")

- Compute new centers for each cluster.

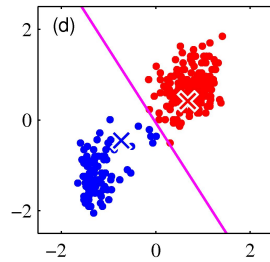


13

## K-Means Clustering

Cluster assignment Step ("E-Step") again

- Re-assign points to the now-nearest center.

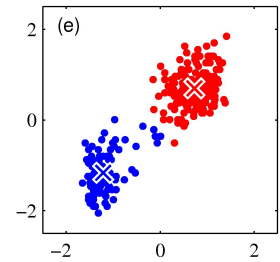


14

## K-Means Clustering

Update parameters (centroids) ("M-Step") again

- Compute centers for the new clusters.

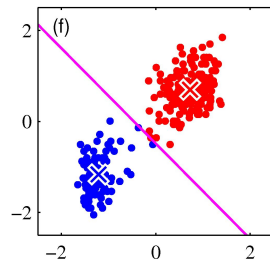


15

## K-Means Clustering

Another Cluster assignment Step ("E-Step")

- Reassign the points to centers.

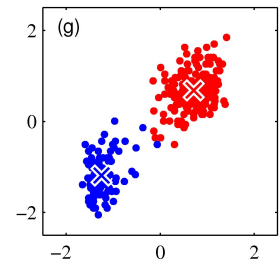


16

## K-Means Clustering

Update parameters (centroids) ("M-Step") again

- New centers.

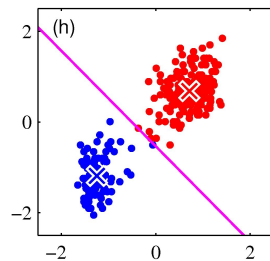


17

## K-Means Clustering

Another Cluster assignment Step ("E-Step")

- New cluster assignments.

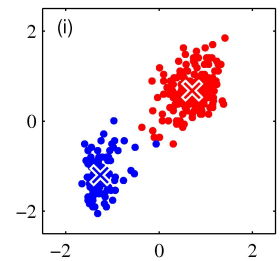


18

## K-Means Clustering

Update parameters (centroids) ("M-Step") again

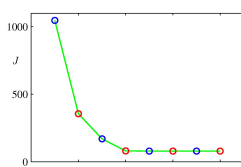
- The cluster centers have stopped changing.



19

## Convergence

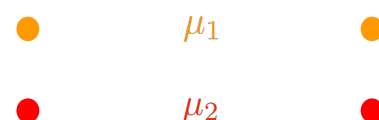
- The objective function of K-means decreases monotonically as the K-means procedure reduces  $J$  in both E-step and M-step.
- Convergence is relatively quick, in steps.
  - blue circles after E-step: assign each point to a cluster
  - red circles after M-step: recompute the cluster centers
  - However, all those distance computations are expensive.



20

## Convergence

- No guarantee that we found the globally optimal solution. The quality of local optimum depends on the initial values.
- The following clustering is a stable local optima



21