

# Section 5: K-means

MCB 112

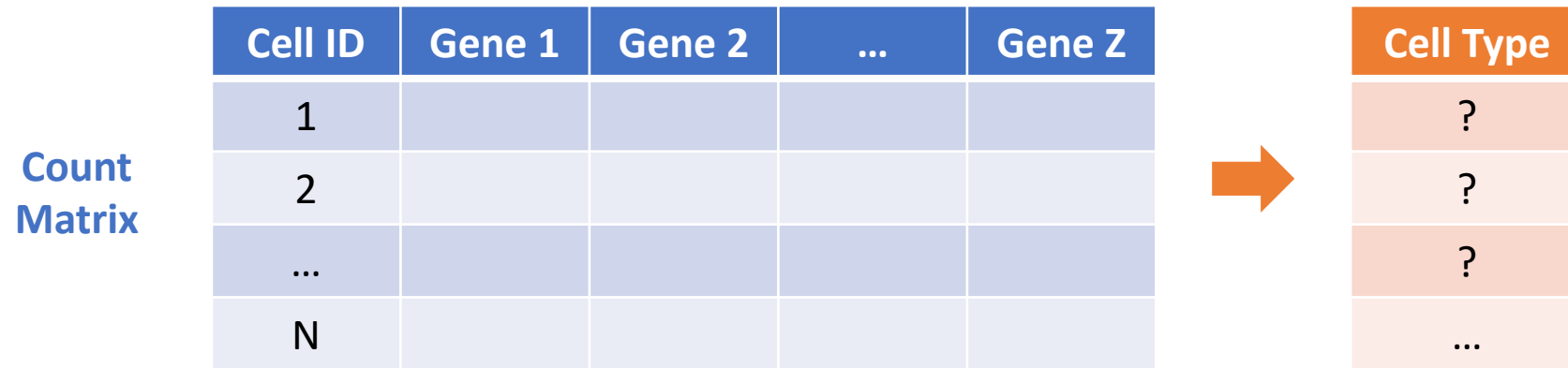
10/07/2022

*(Adapted from 2021 notes)*

# Outline

- Hard and soft K-means
- Mixture Gaussian and Negative Binomial fitting
- Example of 1D hard K-means and 1D mixture gaussian fitting
  - See Jupyter notebook (w05\_section\_jupyter\_notebook.ipynb)

How do we categorize gene counts (i.e. high-dimensional data) into cell type clusters?



## Clustering through K-means/expectation maximization (EM)

### Advantages

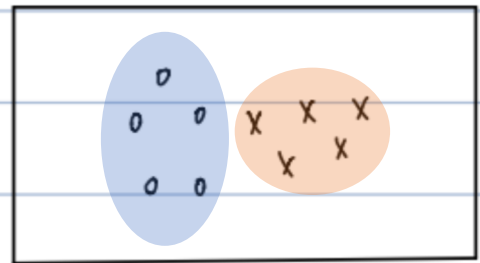
- Simple to implement on a large dataset
- Guarantee convergence
- Generalize to clusters of different shapes and sizes

### Disadvantages

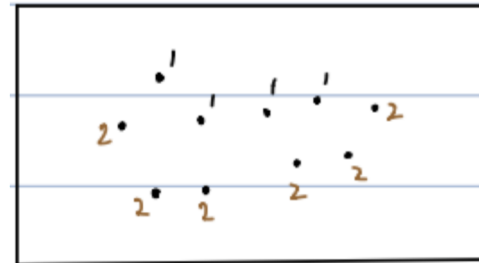
- Choice of K (elbow plot)
- Depend on initial centroid positions (multiple EM runs)
- Influenced by outliers (remove outliers first)

# Simple example of K-means clustering

Ground Truth



Our Data

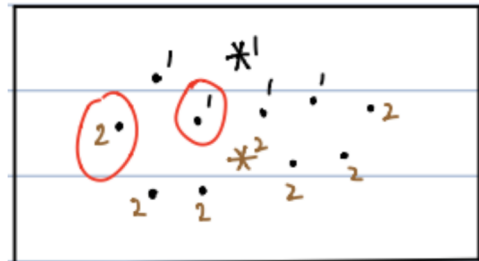
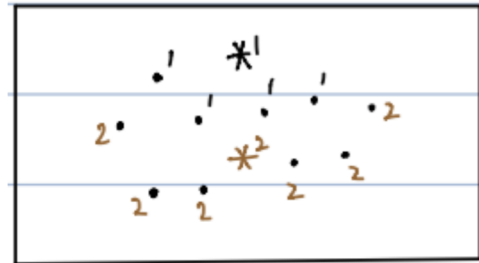


Step 1: Choose  $K$  and

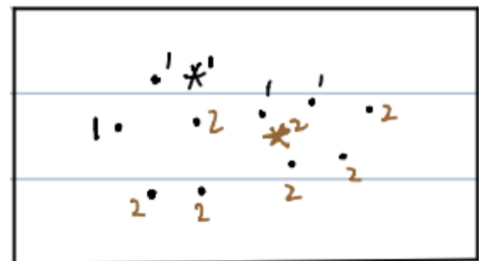
- Assign data randomly to clusters and calculate centroids (shown here)

OR

- Randomly sample centroids and assign data to closest centroids



Step 2: Calculate the distance  $d_{ik}$  between each data point  $i$  and each cluster  $k$ . Update cluster assignment by choosing the closest cluster.



Step 3: Update centroids and iterate until convergence.

# Hard vs. Soft K-means

- Assume we have multi-dimensional data  $\mathbf{X}$ . For each  $X_i$  for  $i = 1, 2, \dots, N$ , it has  $Z$  dimensions (e.g. each cell  $i$  has gene expression counts of  $Z$  genes). We set  $K$  clusters.

## Hard K-means

## Soft K-means

Step 1: Assign data randomly and calculate centroids,  $\mu_{kz}$

$$\mu_{kz} = \frac{\sum_{i \in C_k} X_{iz}}{|C_k|}$$

Step 2: Update cluster assignment based on distances,  $d_{ik}$ , and update centroids,  $\mu_{kz}$

$$d_{ik} = \sqrt{\sum_Z (X_{iz} - \mu_{kz})^2}$$

$$\mu_{kz} = \frac{\sum_{i \in C_k} X_{iz}}{|C_k|}$$

$$r_{ik} = \frac{\exp(-\beta d_{ik}^2)}{\sum_{k'} \exp(-\beta d_{ik'}^2)}$$

where  $\sum_k r_{ik} = 1$  for each  $i$

$$\mu_{kz} = \frac{\sum_i r_{ik} X_{iz}}{\sum_i r_{ik}} \quad (\text{weighted mean of } X_i\text{'s})$$

Step 3: Set a criteria for convergence (objective function)

$$\min \sum_{k=1}^K \sum_{i=1}^N d_{ik}^2$$

Check for convergence

- Same cluster assignment
- Same centroid position (set a threshold)

Step 4: Run multiple times, choose the optimal objective function and corresponding cluster assignment

### Hard K-means

If a point is close to two clusters, we force it to be in one cluster

### Soft K-means

Adjust "stiffness"  $\beta$

- $\beta \rightarrow \infty$ , same as hard K-means
- $\beta \rightarrow 0$ , meaningless

# What to do with an empty cluster?

- Option 1: reinitialize centroid randomly
- Option 2: assign the farthest point as the new centroid

# Mixture model

- $Q$  components (do not need to be the same distribution) with probability  $\pi_q$  for  $q = 1, 2, \dots, Q$
- $\pi_q = P(q | \theta)$  where  $\theta$  is the model parameter(s) and  $\sum_q \pi_q = 1$
- $P(X_i, q | \theta) = P(X_i | q, \theta) \cdot P(q | \theta) = \pi_q \cdot P(X_i | q, \theta)$
- $P(X_i | \theta) = \sum_q P(X_i, q | \theta) = \sum_q \pi_q \cdot P(X_i | q, \theta)$
- $P(\text{cell type} | \text{data } i) = P(q | X_i, \theta) \rightarrow \text{distribution} \left\{ \begin{array}{l} \text{Gaussian} \\ \text{Negative Binomial (RNA-seq)} \end{array} \right.$

# Soft K-means with mixture model

## Soft K-means with mixture model

---

Step 1: Assign data randomly and calculate centroids,  $\mu_q$

Centroids have initial  $\pi_q$

Step 2: Update cluster assignment based on  $P(q | X_i, \theta)$  for each  $q$

$$P(q | X_i, \theta) = \frac{P(X_i, q | \theta)}{P(X_i | \theta)} = \frac{\pi_q \cdot P(X_i | q, \theta)}{\sum_{q'} \pi_{q'} \cdot P(X_i | q', \theta)}$$

- $X_i$  is most likely to be in  $q$  with largest  $P(q | X_i, \theta)$

Step 3: Update  $\mu_q, \pi_q$

$$\mu_q = \frac{\sum_i X_i \cdot P(q | X_i, \theta)}{\sum_i P(q | X_i, \theta)}, \quad \pi_q = \frac{\sum_i P(q | X_i, \theta)}{N}$$



# Soft K-means with mixture model (high-dimensional)

## Soft K-means with mixture model (Z-dimensional)

---

Step 1: Assign data randomly and calculate centroids,  $\mu_{qz}$

Centroids have initial  $\pi_q$

Step 2: Update cluster assignment based on  $P(q | X_i, \theta)$  for each  $q$

$$P(q | X_i, \theta) = \frac{P(X_i, q | \theta)}{P(X_i | \theta)} = \frac{\pi_q \cdot \prod_z P(X_{iz} | q, \theta)}{\sum_{q'} (\pi_{q'} \cdot \prod_z P(X_{iz} | q', \theta))}$$

- $X_i$  is most likely to be in  $q$  with largest  $P(q | X_i, \theta)$

Step 3: Update  $\mu_{qz}, \pi_q$

$$\mu_{qz} = \frac{\sum_i X_{iz} \cdot P(q | X_i, \theta)}{\sum_i P(q | X_i, \theta)}, \quad \pi_q = \frac{\sum_i P(q | X_i, \theta)}{N}$$

# Soft K-means with mixture model

## Soft K-means with mixture model

---

Step 4: Set a criteria for convergence (objective function)

$$\max P(X | \theta) \quad \text{or} \quad \min -\log P(X | \theta)$$

$$P(X | \theta) = \prod_i P(X_i | \theta) \quad (i. i. d)$$

$$= \prod_i \left( \sum_q P(X_i, q | \theta) \right)$$

$$= \prod_i \sum_q \pi_q \cdot P(X_i | q, \theta)$$

$$-\log P(X | \theta) = -\log \left( \prod_i \sum_q \pi_q \cdot P(X_i | q, \theta) \right)$$

$$= -\sum_i \log \left( \sum_q \pi_q \cdot P(X_i | q, \theta) \right)$$



Depend on what distribution we model the data with

Step 5: Run multiple times, choose the optimal cluster assignment

---

# Log-likelihood implementation

$$\begin{aligned} -\log P(X | \theta) &= -\log \left( \prod_i \sum_q \pi_q \cdot P(X_i | q, \theta) \right) \\ &= -\sum_i \log \left( \sum_q \pi_q \cdot P(X_i | q, \theta) \right) \end{aligned}$$



$$\begin{aligned} &\log[\pi_1 P(X_i | \theta_1) + \dots + \pi_Q P(X_i | \theta_Q)] \\ &= \log[e^{\log(\pi_1 P(X_i | \theta_1))} + \dots + e^{\log(\pi_Q P(X_i | \theta_Q))}] \end{aligned}$$



**scipy.special.logsumexp()** :

- Input: an array  $X$
- Compute:  $\log \sum_{x \in X} e^x$

# $P(X_i | q, \theta)$

- Gaussian distribution  $P(X_i | \mu_q, \sigma_q) = \frac{1}{\sqrt{2\pi}\sigma_q} \exp\left(-\frac{(X_i - \mu_q)^2}{2\sigma_q^2}\right)$

- Negative binomial distribution (RNA-seq count data)

- Models the number of failures in a sequence of i.i.d. Bernoulli trials (with success rate  $p$ ) before  $n$  successes occur
- Let  $X$  be the number of failures,  $X \sim \text{NB}(n, p)$

$$P(X = k | n, p) = \binom{k + n - 1}{n - 1} (1 - p)^k p^n$$

- Each cell  $X_i$  has the probability in cluster  $q$  with  $\text{NB}(n, p)$  where

$$n = \frac{1}{\phi}, \quad p = \frac{1}{1 + \mu_q \cdot \phi}$$

models dispersion

updated in each iteration

```
import scipy.stats
scipy.stats.nbinom.logpmf(x, n, p)
```