

Adapted from 2020 notes.

Goals:

1. compare hard and soft kmeans (mixture Gaussian and negative binomial fitting)
2. example of 1D hard kmeans and 1D mixture Gaussian fitting (Jupyter notebook)

Q: How do we categorize/visualize gene counts (high-dimensional data) into cell type clusters?

Example data:

cell id	gene 1 counts	gene 2 counts	...	gene N counts
1				
2				
3				
↓				

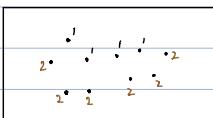
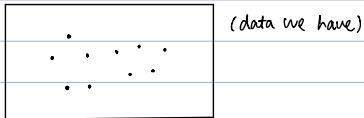
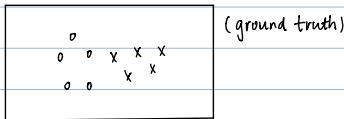
→ cell type?

A: clustering through expectation maximization / kmeans:

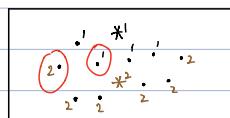
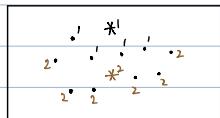
- advantages:
1. simple to implement on large dataset
 2. guarantees convergence
 3. generalize to clusters of different shapes and sizes

- disadvantages:
1. choice of k (loss vs. cluster plot)
 2. depend on initial centroid positions (multiple EM runs)
 3. influenced by outliers (remove first)

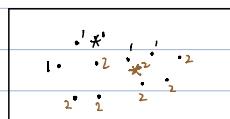
simple example: assume we have two different markers ($0, x$) with 2-D coordinates (x, y)



- ✓
- {
- ① assign data randomly to clusters and calculate centroids
 - ② randomly sample centroids and assign data to closest centroids



point i : cluster 1 2 ... k
distance 3 2 ... d_k



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

Steps:

hard Kmeans ($\geq 1D$)

1D soft Kmeans

1. assign data randomly and calculate centroids

$$\mu_{k,z} = \frac{\sum_{i \in k} X_{i,z}}{|C_k|}$$

once we get d_{ik} , we calculate the responsibility r_{ik} (similar to probability to assign X_i to μ_k)

2. update cluster assignment by calculating distances

$$d_{ik} = \sqrt{\sum_{j=1}^Z (X_{i,j} - \mu_{k,j})^2}$$

$$r_{ik} = \frac{\exp(-\beta d_{ik}^2)}{\sum_{k'=1}^K \exp(-\beta d_{ik'}^2)}$$

update μ_k

$$\mu_{k,z} = \frac{\sum_{i \in k} X_{i,z}}{|C_k|}$$

where $\sum_{k'=1}^K r_{ik} = 1$ for each i

3. set a criteria for convergence (objective function)

$$\sum_{k=1}^K \sum_{i=1}^N d_{ik}^2 \text{ (minimize)}$$

$$\mu_k = \frac{\sum_{i=1}^N r_{ik} \cdot X_i}{\sum_{i=1}^N r_{ik}} \text{ (weighted average of all } X_i)$$

1) check if point assignments change

2) check if centroid locations change

local minimizer
need to set a reasonable threshold

3) check if objective function changes

4. run Kmeans multiple times and choose optimal objective function and corresponding cluster assignment

hard: if a point is close to 2 clusters, we force it to be only one cluster

stiffness β $\begin{cases} \rightarrow \infty : \text{hard Kmeans} \\ \rightarrow 0 : \text{useless} \end{cases}$

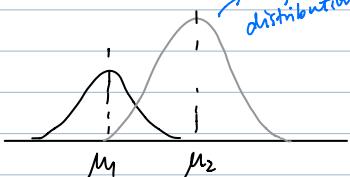
① reinitialize randomly

Also mention the case of empty cluster: ② find the farthest point and assign it as the new centroid

Mixture model:

can be

any distribution



1. Q components with probability π_q for $q=1, 2, \dots, Q$
model parameters

$$\pi_q = P(q|\theta) \text{ and } \sum_{q=1}^Q \pi_q = 1$$

$$2. P(X_i, q|\theta) = P(X_i|q, \theta) \cdot P(q|\theta) = \pi_q \cdot P(X_i|q, \theta)$$

$$P(X_i|\theta) = \sum_{q=1}^Q P(X_i, q|\theta) = \sum_{q=1}^Q \pi_q P(X_i|q, \theta)$$

Q: What's the probability of cell i belonging to certain cell type given the RNA seq data?

$$P(\text{cell type} | \text{data } i) = P(q | X_i, \theta)$$

↑ distribution { Gaussian
Negative binomial (RNA seq)}

Soft kmeans with mixture model

Steps:

1. assign data randomly and calculate centroids

Centroids have initial π_q and μ_q

depend on what distribution we will use

$$P(q | X_i, \theta) = \frac{P(X_i, q | \theta)}{P(X_i | \theta)} = \frac{\pi_q \cdot P(X_i | q, \theta)}{\sum_{q=1}^Q \pi_q \cdot P(X_i | q, \theta)}$$

2. update cluster assignment by calculating $P(q | X_i, \theta)$

for each q .

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

update μ_q

$$\hat{\mu}_q = \frac{\sum_{i=1}^N X_i \cdot P(q | X_i, \theta)}{\sum_{i=1}^N P(q | X_i, \theta)} \quad \text{and} \quad \hat{\pi}_q = \frac{\sum_{i=1}^N P(q | X_i, \theta)}{N}$$

3. set a criteria for convergence (objective function)

We want to maximize $P(X | \theta)$ or minimize $-\log P(X | \theta)$

$$P(X | \theta) = \prod_{i=1}^N P(X_i | \theta) \quad (\text{i.i.d})$$

$$= \prod_{i=1}^N \left[\sum_{q=1}^Q P(X_i, q | \theta) \right]$$

$$= \prod_{i=1}^N \sum_{q=1}^Q \pi_q \cdot P(X_i | q, \theta)$$

$$-\log P(X | \theta) = -\log \left(\prod_{i=1}^N \sum_{q=1}^Q \pi_q \cdot P(X_i | q, \theta) \right)$$

$$= -\sum_{i=1}^N \log \left(\sum_{q=1}^Q \pi_q \cdot P(X_i | q, \theta) \right)$$

logsumexp function: $\log \sum_{x=x} e^x$

its input is an array of x

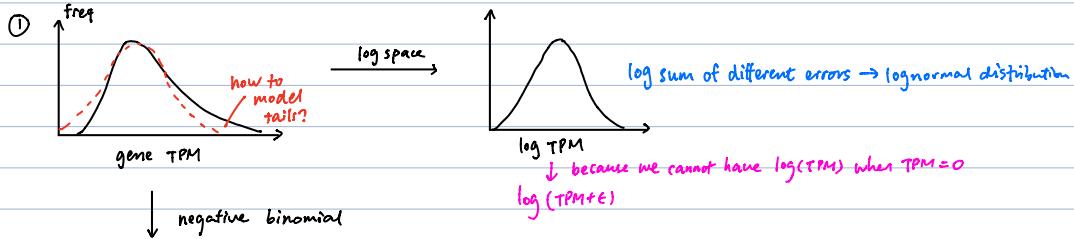
$$\begin{aligned} & \log \left[\pi_1 P(X_i | \theta_1) + \pi_2 P(X_i | \theta_2) + \dots + \pi_Q P(X_i | \theta_Q) \right] \\ &= \log \left[e^{\log(\pi_1 \cdot P(X_i | \theta_1))} + \dots + e^{\log(\pi_Q \cdot P(X_i | \theta_Q))} \right] \end{aligned}$$

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

assume constant, but may not be

$$\text{if Gaussian: } P(X_i | \mu_q, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(X_i - \mu_q)^2}{2\sigma^2}\right)$$

For scRNA-seq, we can model counts with either lognormal / negative binomial distributions.



② Poisson $\sim (\mu)$ → it's discrete while TPM is continuous. (NB is better in modeling counts)

+ sampling a transcript from a large pool has a Poisson distribution
NB to model the dispersion and the uncertainty in the mean of the distribution,
↓ we use gamma-poisson distribution with 2 parameters

gamma-poisson $\sim (\mu, \phi)$

↓
mean # of occurrence
of a rare event

Define NB distribution: # of failures in a series of i.i.d Bernoulli trials before n successes occurs, with success rate of p.
let X be the # of failures, $X \sim NB(n, p)$

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

Each gene X_i has prob in cluster q with $NB(n, p)$ where $n = \frac{1}{\phi}$ and $p = \frac{1}{1 + \mu_q \cdot \phi}$
↑ models dispersion updated

$$\begin{aligned} &\text{in } P(X_i | q, \theta) \\ &= P(X_i | \mu_q) \end{aligned}$$

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