

Mixture Model

DATASET : $\{x_i\}_{i=1}^N$

DATA Likelihood :

$$P(x|\theta) = \sum_{k=1}^K \pi_k \cdot P(x|\theta_k)$$

$K = \# \text{ components}$ (hyperparameter)

where $0 \leq \pi_k \leq 1 \quad \forall k$

and $\sum_{k=1}^K \pi_k = 1$

$\theta = \text{set of all parameters for all } K$
components

$$\theta = \{\theta_k\}_{k=1}^K$$

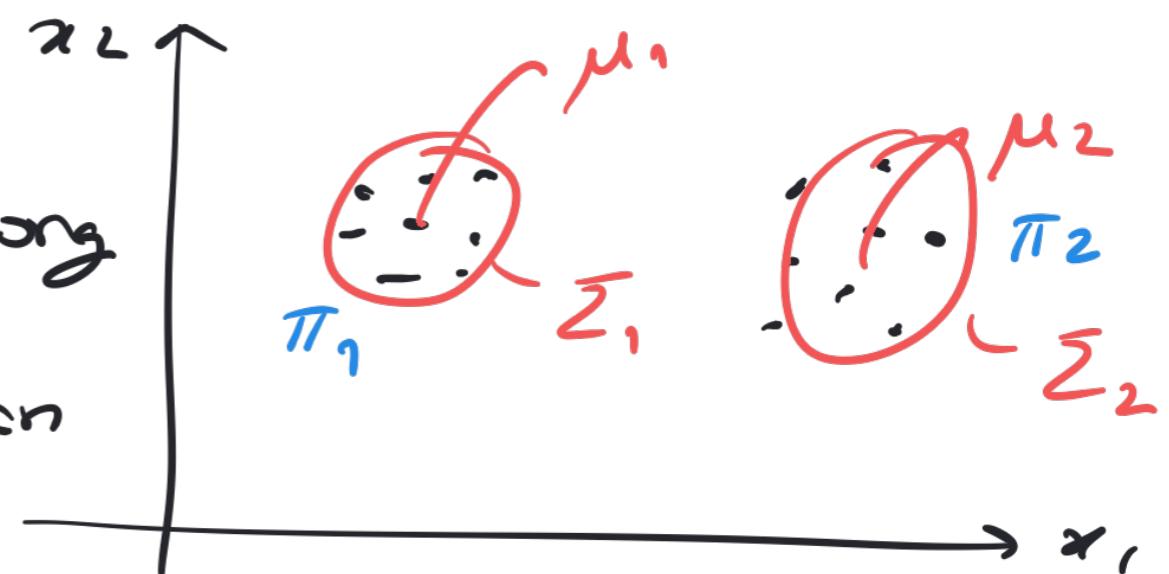
Gaussian Mixture Models (GMMs)

$$P(x|\theta) = \sum_{k=1}^K \pi_k \cdot N(x|\mu_k, \Sigma_k)$$

where $0 \leq \pi_k \leq 1$ and $\sum_{k=1}^K \pi_k = 1$

PARAMETERS: $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

each x_i can only belong
to a single Gaussian
component.



$$\text{DATASET} = \{x_i\}_{i=1}^N$$

OBSERVED DATA (i.i.d.): $\mathcal{L}^o = \prod_{i=1}^N P(x_i | \theta)$
Likelihood

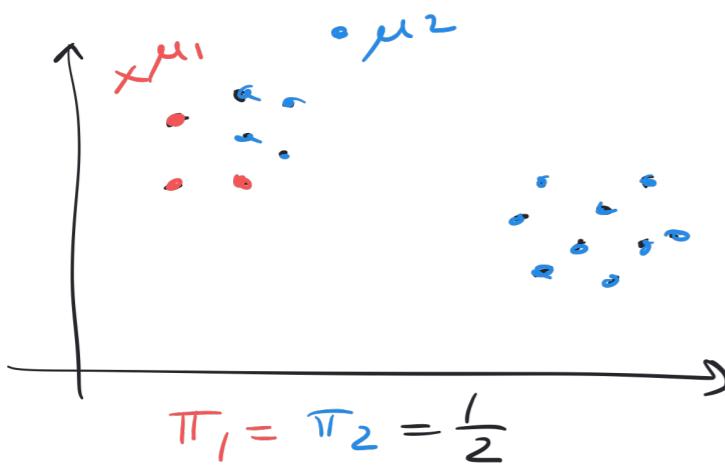
$$\mathcal{L}^o = \prod_{i=1}^N \sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k)$$

STEP 1: Provide a value for # components

$$\text{e.g. } K=2$$

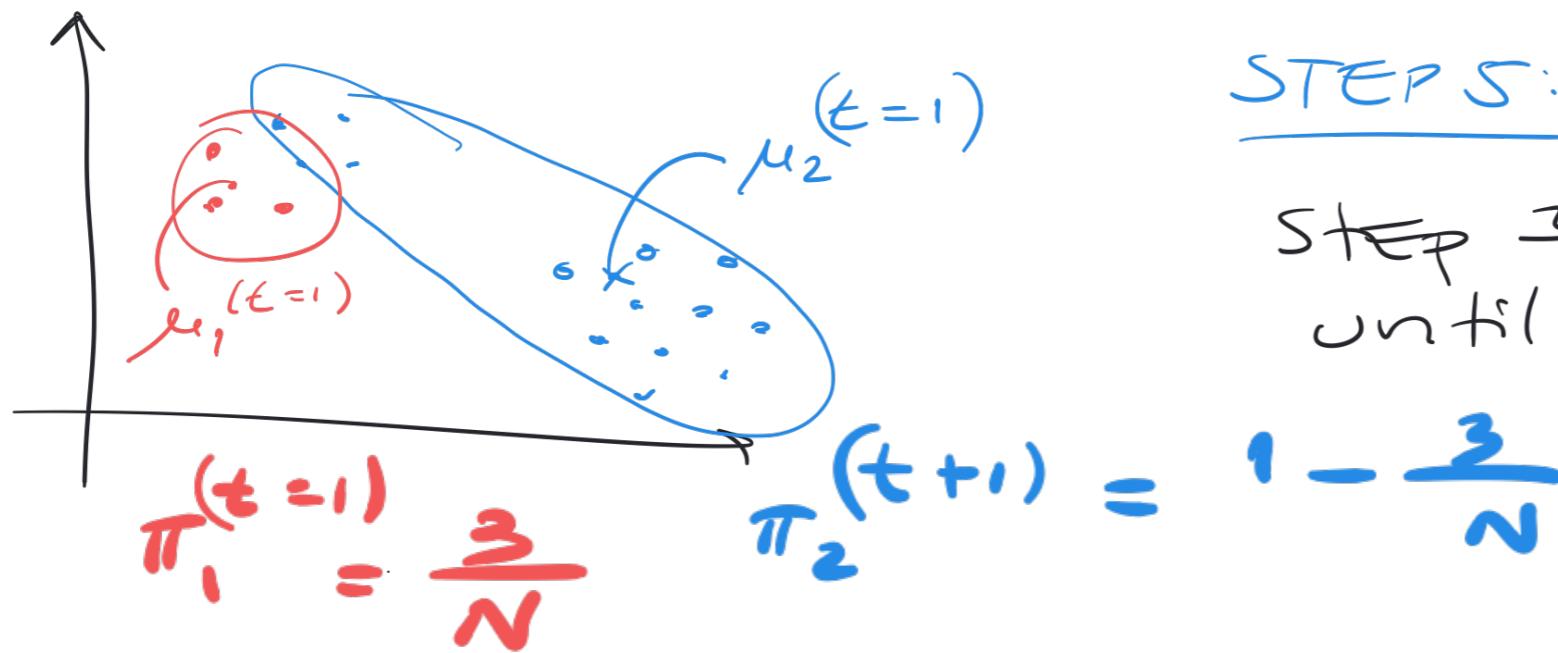
STEP 2: Initialize the parameters $\Theta^{(t)} = \{\pi_k, \mu_k, \Sigma_k\}_K$

$t \equiv \text{iteration}$



STEP 3: Assign a membership value for every point.

STEP 4: Fix membership value and update the parameters Θ .



STEP 5: Go back to STEP 3 and continue until convergence

$$\frac{1}{\sqrt{2\pi} |\Sigma_k|^{1/2}} N(x | \mu_k, \Sigma_k)$$

Gaussian component

Applications of GMM:

① DATA

like likelihood

Estimation.

② clustering.

Partition the data into clusters
based on density similarity.

In MAP, we can include a prior
on any of the parameters

$$\boldsymbol{\theta} = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K \quad \text{but it}$$

will not have any conjugate

prior relationship.

Observed DATA likelihood:

$$\mathcal{L}^o = \prod_{i=1}^N \sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k)$$

$$\begin{aligned}\mathcal{L} &= \ln(\mathcal{L}^o) \\ &= \sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k) \right)\end{aligned}$$

this is a "difficult" problem.

$$\frac{\partial \mathcal{L}}{\partial \mu_k} = 0$$

Expectation - Maximization (EM) Algorithm

- ① Introduce Hidden Latent Variables that will simplify the problem.
- ② Iterative greedy algorithm
 - ↳ Solution will depend on the initialization
 - ↳ it may converge to a local optima.

Example : Censored data

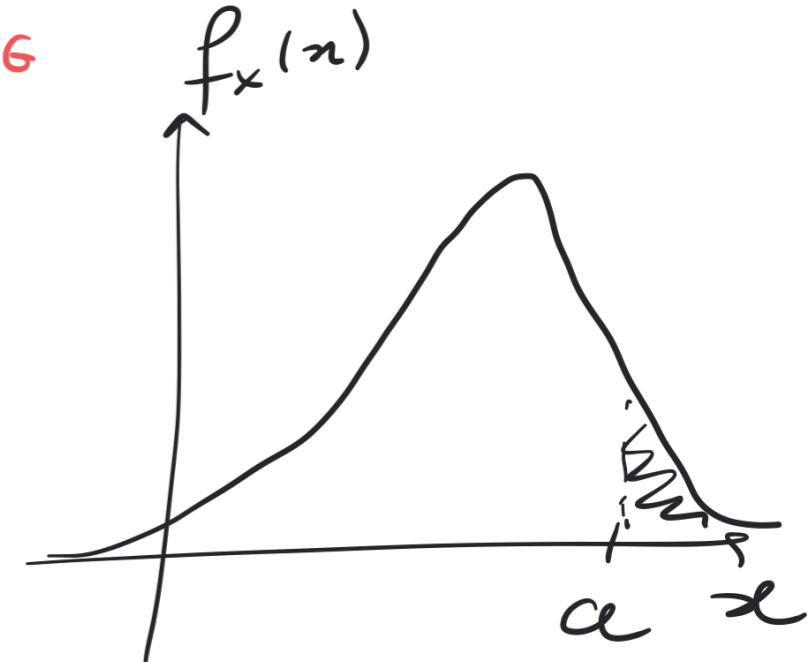
$$\underbrace{x_1, x_2, \dots, x_m}_{\equiv \text{measure the true value}}$$

\equiv measure the true value

$$\underbrace{x_{m+1}, \dots, x_N}_{=a = \text{censored}}$$

at a fixed value

TRAINING DATA



Prob.
mass
fct

$$P_x(x) = \begin{cases} x, & x < a \\ a, & x \geq a \end{cases}$$

OBSERVED

DATA

Likelihood

$$\mathcal{L}^o = \prod_{i=1}^m P(x_i|\theta) \cdot \prod_{j=m+1}^N \left(\sum_{x_j=a} P(x_j|\theta) \right)$$

EM introduces the hidden latent variables z that, if given, would simplify the problem.

In the censored data example:

z_j = true (unknown) values for censored samples, $j = m+1, \dots, N$

with this \geq variable:

$$\mathcal{L}^0 = \prod_{i=1}^m P(x_i | \theta) \cdot \prod_{j=m+1}^N P(z_j | \theta)$$

$$[\mathcal{L} = \ln \mathcal{L}^0]$$

To find z 's, the EM algorithm will
find the expectation value (E-STEP)
for data likelihood given some
initialization for parameters θ .

Optimization function to search for

z 's given initial values for

parameters $\theta^{(t)}$ ($t \rightarrow$ iteration) :

expectation

$$Q(\theta, \theta^{(t)}) = E_z [\ln(L^0) | X, \theta^{(t)}]$$

SET OF
UNKNOWN
PARAMETERS

initializations
(numerical
values)

X = training
data or
feature
matrix

///
 X is a discrete R.V.

$\mathbb{E}_x[X] = \text{mean of } X$

$$\mathbb{E}_x[X] = \sum_x x \cdot p_x(x)$$

$$\mathbb{E}_x[f(x)] = \sum_x f(x) \cdot p_x(x)$$

///

$$Q(\theta, \theta^t) = E_z [\ln \mathcal{L}^\circ | x, \theta^t]$$

$$Q(\theta, \theta^t) = \sum_{z_i} \ln(\mathcal{L}^\circ) \cdot P(z_i | x, \theta^t)$$

EM

$$\arg \max_{\{\theta, z\}} Q(\theta, \theta^t)$$

Pseudo-code for EM algorithm

- ① Initialization $t = 0$
- ② Initialize parameters θ , $\theta^{(t)} = \{\mu_k^{(t)}, \Sigma_k^{(t)}, \pi_k^{(t)}\}$
- ③ E-STEP: Hold $\theta^{(t)}$ fixed and find
 z' 's that maximize $Q(\theta, \theta^{(t)})$
- ④ M-STEP (Maximization step): Hold z' 's
for step ③ and find new values for
 θ .
- ⑤ Iterate between steps ③ and ④ until
convergence.

convergence criteria:

- (S.1) Maximum # iterations
- (S.2) threshold for the difference in
parameter values between 2 consecutive
iterations.

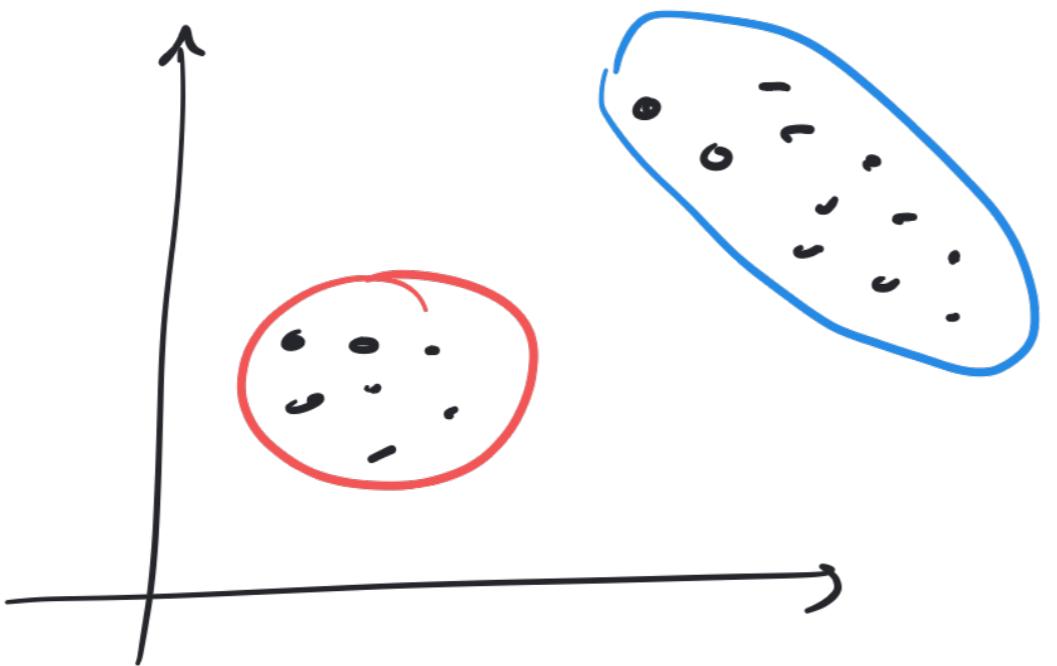
EM algorithm is an

Alternating optimization

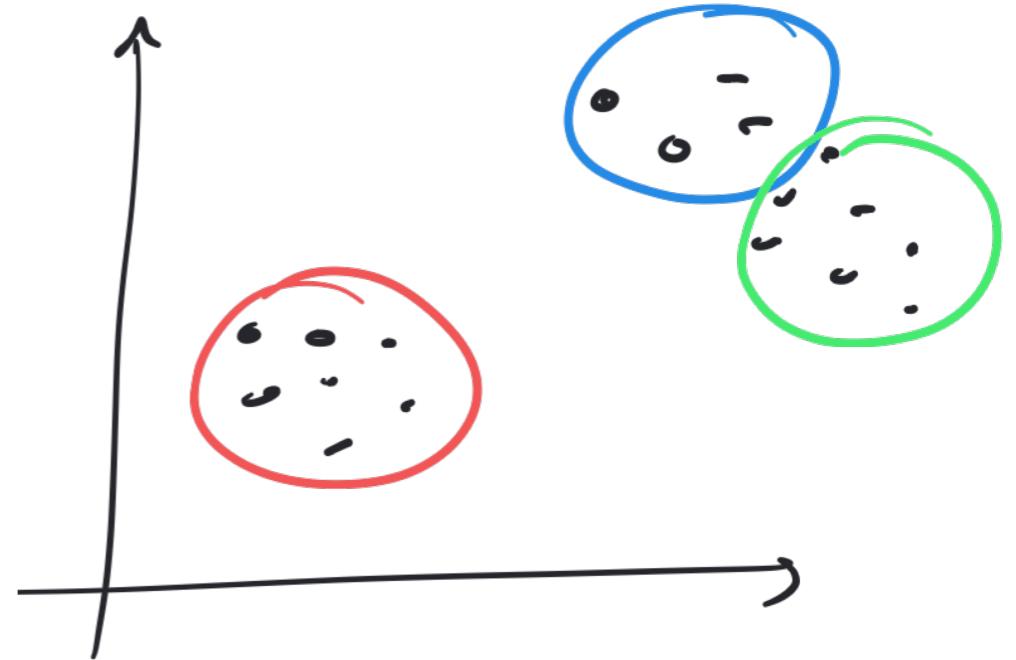
algorithm.

① sensitive to the initialization of parameters θ .

② Convergence is not guaranteed.



$K = 2$



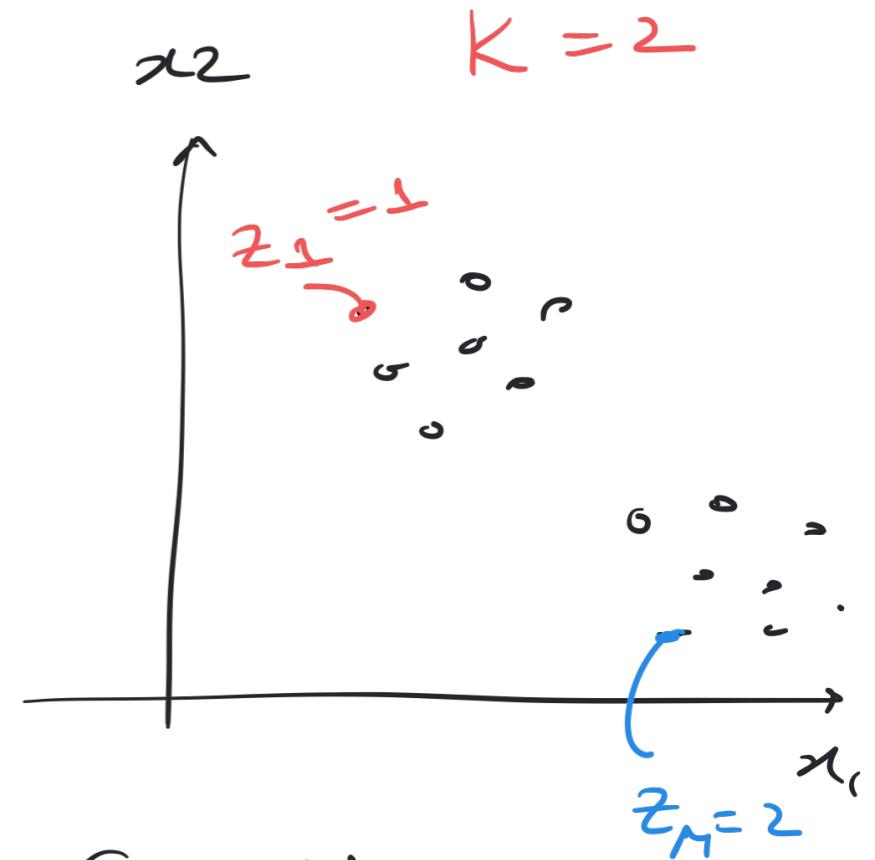
$K = 3$

GMM

$$\mathcal{L}^0 = \prod_{i=1}^N \sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k)$$

where $0 \leq \pi_k \leq 1$ and $\sum_k \pi_k = 1$

$$\text{and } \theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$$



Hidden Latent variables: $z_i \equiv$ label of Gaussian component from which point x_i was drawn from

$$z_i \in \{1, 2, \dots, K\}$$

Complete DATA Likelihood:
(this includes z 's)

$$\mathcal{L}^c = \prod_{i=1}^N \pi_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

EM Optimization:

$$Q(\theta, \theta^{(t)}) = E_z [\ln \mathcal{L}^c | X, \theta^{(t)}]$$

$$Q(\theta, \theta^t) = \mathbb{E}_z [\ln \mathcal{L}^c | x, \theta^t]$$

$$= \sum_{z_i=1}^K \ln(\mathcal{L}^c) \cdot P(z_i | x, \theta^t)$$

$$= \sum_{z_i=1}^K \ln \left(\prod_{i=1}^N \pi_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i}) \right).$$

E-STEP: Fix the parameters so this term is constant

$$P(z_i | x, \theta^t)$$

M-STEP: this term will be constant.

$$\mathcal{L}^c = \prod_{i=1}^N \pi_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

where

$$N(x_i | \mu_{z_i}, \Sigma_{z_i}) = \frac{1}{(2\pi)^{d/2} |\Sigma_{z_i}|^{1/2}} \exp\left(-\frac{1}{2} (x_i - \mu_{z_i})^\top \Sigma_{z_i}^{-1} (x_i - \mu_{z_i})\right)$$

$d \equiv$ dimension of feature space.

$$Q(\theta, \theta^t) = \sum_{\substack{z_i=1 \\ z_i=1}}^K \ln \left(\prod_{i=1}^N \pi_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i}) \right) \cdot P(z_i | x, \theta^t)$$

$$= \sum_{k=1}^K \ln \left(\prod_{i=1}^N \pi_k \cdot N(x_i | \mu_k, \Sigma_k) \right) \cdot P(z_i = k | x, \theta^t)$$

$$= \sum_{k=1}^K \left[\sum_{i=1}^N \left(\ln(\pi_k) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right) \right] \cdot P(z_i = k | x, \theta^t)$$

Let's assume $\Sigma_k = \sigma_k^2 \cdot H$

$$|\Sigma_k| = (\sigma_k^2)^d$$

$$\Sigma_k^{-1} = \frac{1}{\sigma_k^2} \cdot H$$

Plugging in:

$$Q(\theta, \sigma^2) = \sum_{k=1}^K \left[\sum_{i=1}^n \left[\ln(\pi_k) - \frac{d}{2} \ln(2\pi) - \frac{\alpha}{2} \ln(\sigma_k^2) - \frac{1}{2\sigma_k^2} \|x_i - \mu_k\|^2 \right] \right] \cdot P(z_i=k|x, \theta^t)$$

① Initialize the parameters

$$\theta^{(t)} = \left\{ \pi_k^{(t)}, \mu_k^{(t)}, \sigma_k^{(t)} \right\}_{k=1}^K$$

② E-step

Hold $\theta^{(t)}$ fixed and solve for

$$P(z_i=k|x, \theta^t) = c_{ik}$$

= membership
of sample x_i
in component k

$$C = \begin{bmatrix} 1 & 0.7 & 0.1 & \dots & 0.2 \\ 2 & : & : & & : \\ \vdots & & & & \vdots \\ n & & & & \end{bmatrix} \xrightarrow{\sum_{k=1}^K C_{ik} = 1}$$

Soft membership assignment.

E-Step

$$P(z_i | x_i, \theta^{(t)}) = \frac{P(x_i | z_i, \theta^{(t)}) \cdot P(z_i | \theta^{(t)})}{P(x_i | \theta^{(t)})}$$

$$= \frac{P(x_i | z_i, \theta^{(t)}) \cdot P(z_i | \theta^{(t)})}{\sum_{k=1}^K P(x_i | z_i = k, \theta^{(t)}) \cdot P(z_i = k | \theta^{(t)})}$$

$$= \frac{N(x_i | \mu_{z_i}^{(t)}, \Sigma_{z_i}^{(t)}) \cdot \pi_{z_i}^{(t)}}{\sum_{k=1}^K N(x_i | \mu_{z_k}^{(t)}, \Sigma_k^{(t)}) \cdot \pi_k^{(t)}} = c_{ik}$$