

Lecture 6/13: Generative Models so far

train $P(x)$ = density / curve [param] fit to data x (Sep for each Y)

predict $P(Y|x) = \frac{P(x|Y) \cdot P(Y)}{P(x)}$

• For each Y $P(x|Y) = \mathcal{N}(x | \mu, \Sigma)$ ^{D-dim} \Rightarrow Gauss DA

• For each Y $P(x|Y) = P(x^1 x^2 \dots x^D | Y)$
 $= P(x^1|Y) \cdot P(x^2|Y) \dots P(x^D|Y) \Rightarrow$ Naive Bayes

• Not-So-Naive: \rightarrow **Belief Network** group features that are very dependent

Bayes \Rightarrow **indep**

For each Y : $P(x|Y) = P(x^1 x^2 \dots x^D | Y) =$

$= P(x^1, x^2 | Y)$ \cdot $P(x^3 | Y)$ \cdot $P(x^4, x^5 | Y, x^1)$ \cdot $P(x^6, x^7 | Y)$

2dim *1-dim* *2dim fit for Y, x^1*

MIXT OF GAUSSIANS (GMM)

Today +

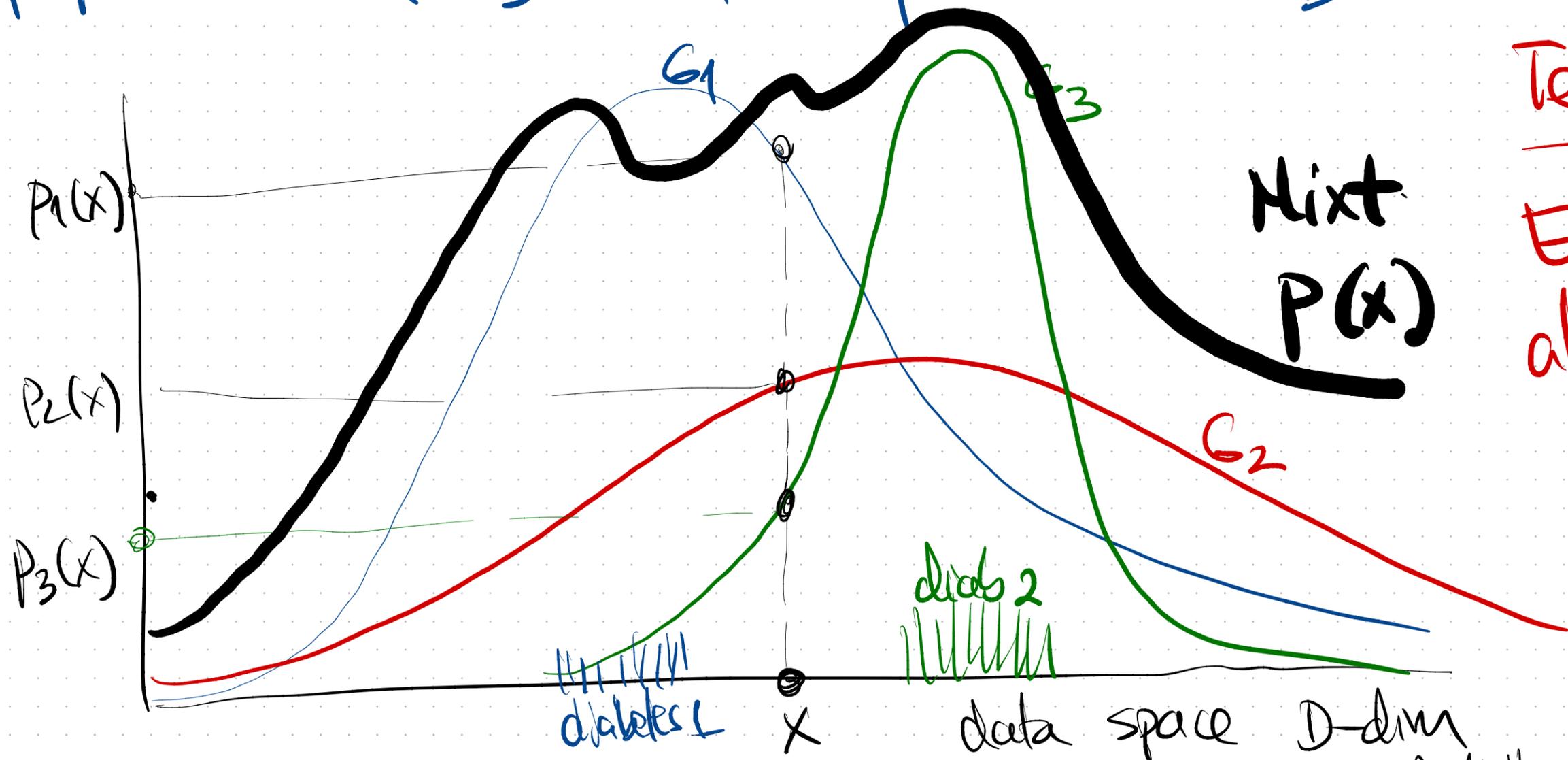
Next:

$$P(x) = \underbrace{w_1}_{\text{fixed}} \underbrace{N_1(x | \mu_1, \Sigma_1)}_{D\text{-dim}} + \underbrace{w_2}_{\text{fixed}} \underbrace{N_2(x | \mu_2, \Sigma_2)}_{D\text{-dim}} + \underbrace{w_3}_{\text{fixed}} \underbrace{N_3(x | \mu_3, \Sigma_3)}_{D\text{-dim}}$$

$w_1 + w_2 + w_3 = 1$
proportions

$K=3$ # of "components"

$D = \# \text{ data dim}$



Technique

$E(x, p)$ Max
algorithm

Mixture adv: ability to fit data with multiple "hills"
"bumps"

Recap for intuition: Kmeans clustering algorithm.

X is data \Rightarrow cluster into $K=6$ groups
 $i=1:N$ $k=1:6$

π_{ik} = membership
 $\begin{cases} 1 & \text{if } x_i \rightarrow \text{cluster } k \\ 0 & \text{if not} \end{cases}$

E-step
 calculate π_{ik} , given μ_k

M-Step
 calculate μ_k , given π_{ik}

μ_k = centroid for cluster k
param

$$\pi_{ik} = 1 \text{ for closest } \mu_k$$

$$= \underset{k}{\text{argmin}} \|x_i - \mu_k\|^2$$

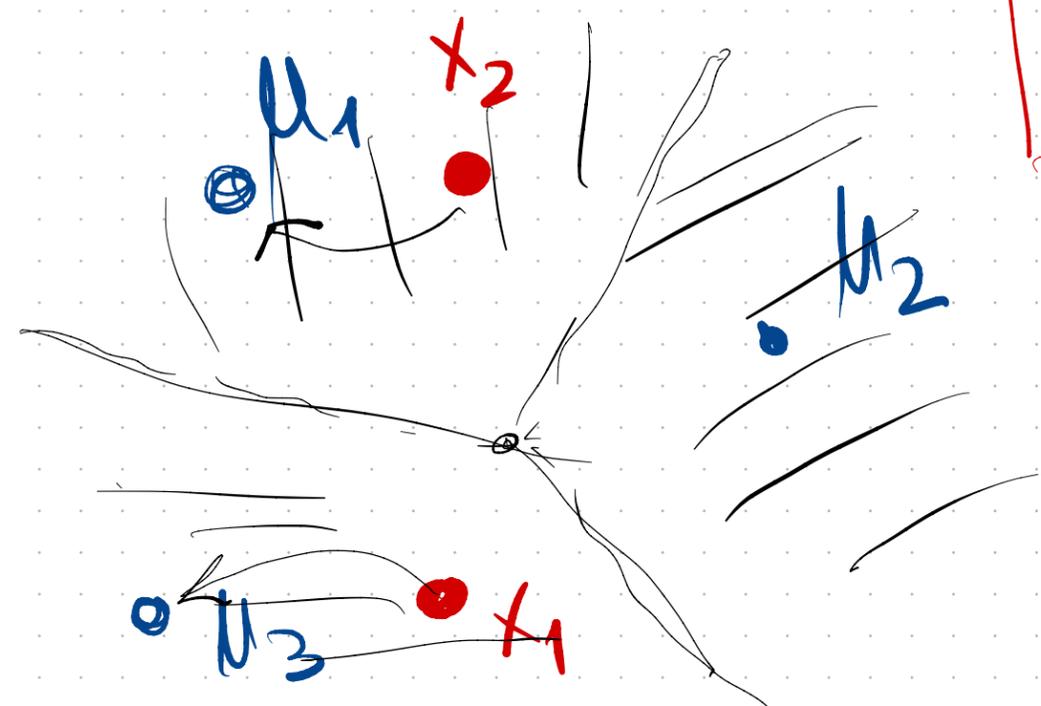
$$\mu_k = \text{avg of } (x \text{ points}) \text{ in cluster } k$$

$$= \frac{\sum_{i=1}^N x_i \cdot \pi_{ik}}{\sum_{i=1}^N \pi_{ik}}$$

filter out points in cluster k

$$\sum_{i=1}^N \pi_{ik} = N_k = \# \text{ of points in cluster } k$$

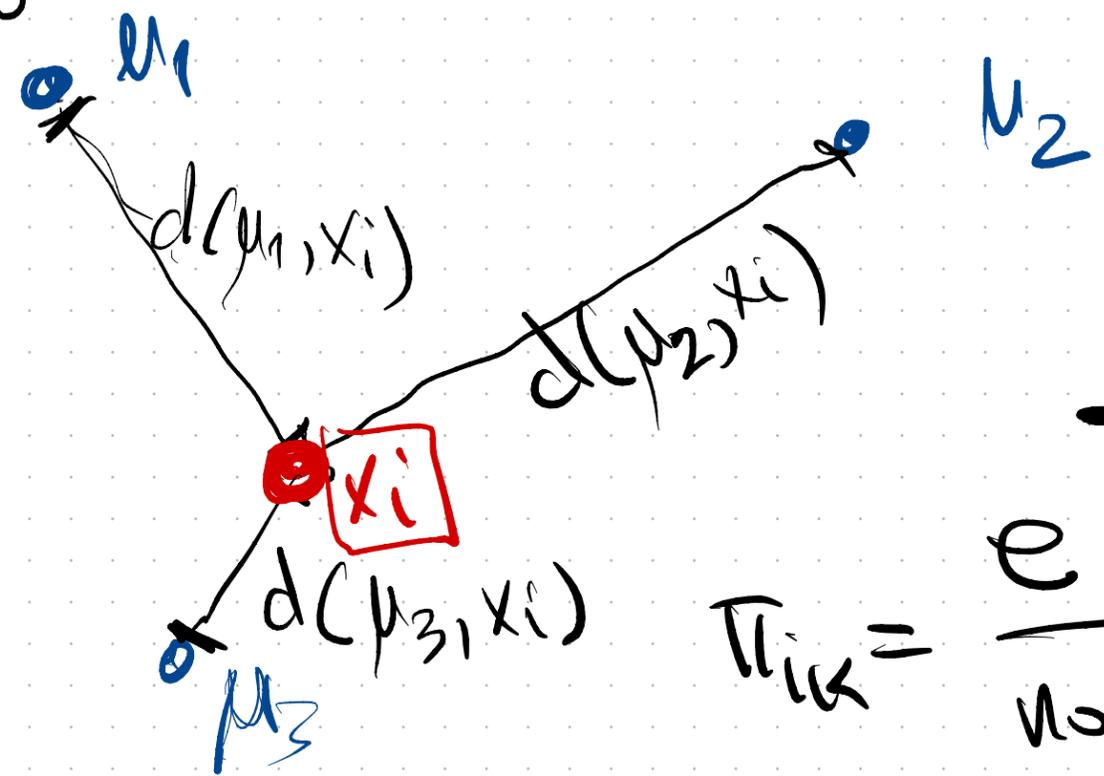
HARD: $\pi_{ik} \rightarrow \mu_k$



SFT: $\pi_{ik} = \text{prob}(x_i \rightarrow \mu_k)$

E step

- score dist \Rightarrow probability
- reg output



M step (same)

$$\mu_k = \frac{\sum_i x_i \cdot \pi_{ik}}{\sum_i \pi_{ik}}$$

proportion (weight)

\rightarrow weighted avg

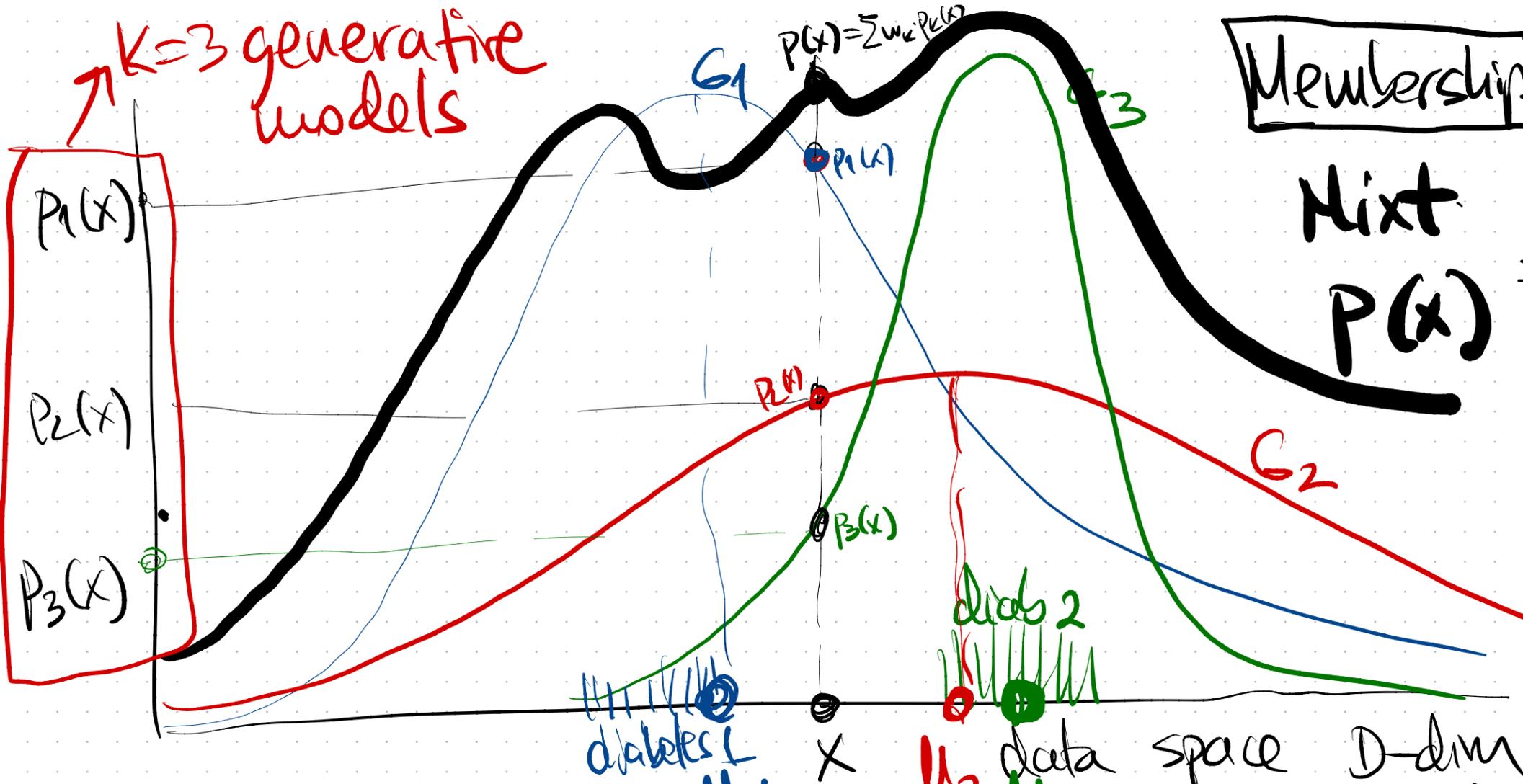
$N_k = E[\# \text{ points}]$ in cluster k

$$\pi_{ik} = \frac{e^{-\frac{1}{\beta} \|x_i - \mu_k\|^2}}{\text{normalized over all } k}$$

fix datapoint x_i small dist $\|x_i - \mu_k\|^2 \equiv$ large π_{ik}

$[\pi_{i1} \ \pi_{i2} \ \dots \ \pi_{ik}]$ distribution over clusters

- parts of x_i distributed $\rightarrow x_i \cdot \pi_{ik}$
- which cluster/group x_i is in \rightarrow which generator μ_k produced datap x_i



$k=3$ generative models

Membership π_{ik} = prob that x_i was generated by Gauss k

Mixt $P(x) = \sum_{k=1}^K w_k \cdot P(x_i | \mu_k, \Sigma_k)$

Normalized over all k

CONST: $w_k =$ size of cluster / source k

$E[\# \text{datap}]$ from k

$= \sum_{i=1}^N \pi_{ik}$

Mixture adv: ability to fit data with multiple "hills" μ_{clusters}

$$P(x) = \sum_{k=1}^3 w_k \cdot \mathcal{N}_k(x | \mu_k, \Sigma_k)$$

D -dim

prob to observe patient/email/car x_i overall sources $k=1, k=2, \dots, k=K$

π_{ik} = probab of source/generator k for observed datap x_i (model the evidence)
 w_k = prob of source/gen k in general (prior) \rightarrow given x_i

E step

calculate π_{ik} given
 params $(\mu_k, \Sigma_k, w_k)_{k=1:K}$

$$\pi_{ik} = \text{pr}(k \text{ source generated } x_i)$$

$$= \text{pr}(k | x_i) = \frac{\text{pr}(x_i | k) \cdot \text{pr}(k)}{\text{pr}(x_i)}$$

$$= \frac{N_k(x_i | \mu_k, \Sigma_k) \cdot w_k}{\text{normalize over } k \text{ generators}}$$

normalize over k
 generators

M step

calculate param $(\mu_k, \Sigma_k, w_k)_{k=1:K}$
 given memberships π_{ik}

- Loglikelihood (data) = $LL(x)$

- take expectation w.r.t π_{ik}
 $E[LL(x)] \approx LL(x)$

- take differentials Constraint $w_1 + w_2 + w_3 = 1$

$$\frac{\partial LL(x)}{\partial \mu_k} = 0 \quad \left| \quad \frac{\partial LL(x)}{\partial \Sigma_k} = 0 \right| \quad \frac{\partial LL(x)}{\partial w_k} = ?$$

Lagrangean