

Lecture 6/13 : Generative Models so far

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

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

• For each c $P(x|c) = w(x|\mu, \Sigma)$ $\xrightarrow{D\text{-dim}}$ Gauss DA

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

• Not-So-Naive: $\xrightarrow{\text{Belief Network}}$ group features that are very dependent
Bayes

$$\begin{aligned} \text{For each } c: P(x|c) &= P(x^1|x^2, \dots, x^D|c) \\ &= P(x^1, x^2|c) \cdot P(x^3|c) \cdot P(x^4, x^5|x^1) \cdot P(x^6, x^7|c) \end{aligned}$$

indep

2dim

1-dim

2dim fit for x_2, x_3

MIXT OF GAUSSIANS (GMM)

Today:

$$\text{Next: } P(x) = w_1 N_1(x | \mu_1, \Sigma_1) + w_2 N_2(x | \mu_2, \Sigma_2) + w_3 N_3(x | \mu_3, \Sigma_3)$$

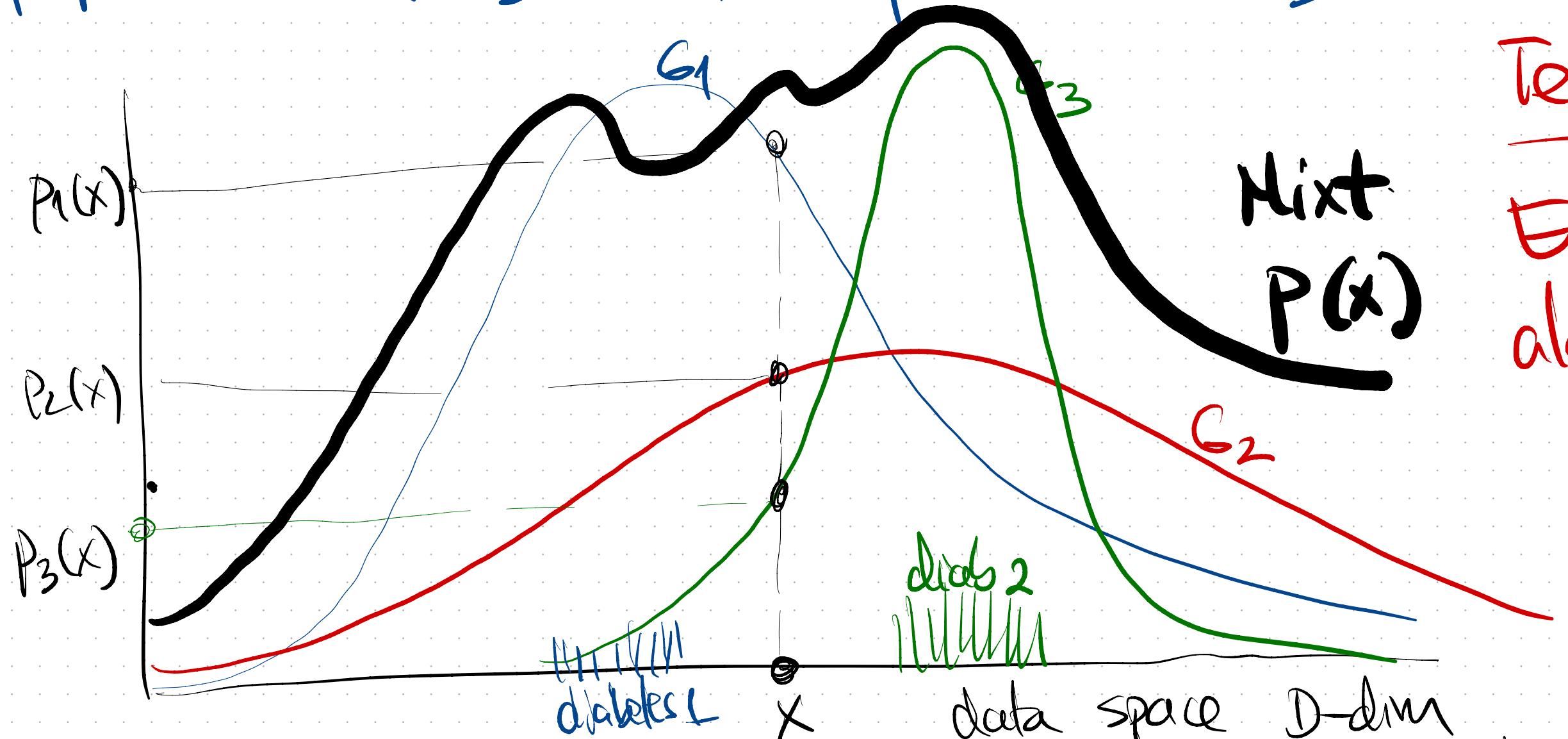
w₁ w₂ w₃ = 1 fixed D-dim fixed D-dim fixed D-dim

$$w_1 + w_2 + w_3 = 1$$

proportions

K=3 # of "components"

D = # data dim



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

$$\text{GMM} = w_1 \cdot w_1() + w_2 \cdot w_2() + w_3 \cdot w_3() \quad k=3$$

- "generative process": assume that datapoint x_i is generated indep of all other datapoints, in 2 steps.

1) select one comp (Gauss) with fix probab $w_1/w_2/w_3$
generator $N_k(\mu_k, \Sigma_k)$

2) sample x_i from N_k density (D-dim)

- Think of K-means (hard) and soft-kmeans mechanics

• π_{ik} = membership

R.V. $\begin{cases} 1 & \text{if } x_i \text{ generated by } W_k() \\ 0 & \text{if not} \end{cases}$

prob $\pi_{ik} = \text{prob that } W_k() \text{ generated } x_i$
 $\pi_{ik} = \langle \pi_{ik} \rangle = E[\pi_{ik}] = \text{prob}(\pi_{ik}=1)$

fixed i: $\sum_{k=1}^K \langle \pi_{ik} \rangle = 1$

fixed k (cluster/generator)

$$\sum_{i=1}^N \langle \pi_{ik} \rangle = E[\# \text{ datapoints generated by } W_k]$$

Recap for intuition: Kmeans clustering algorithm.

X D-dim data \Rightarrow clusters mb $K=6$ groups
 $i=1 \dots N$ $k=1 \dots 6$

- 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

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

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

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

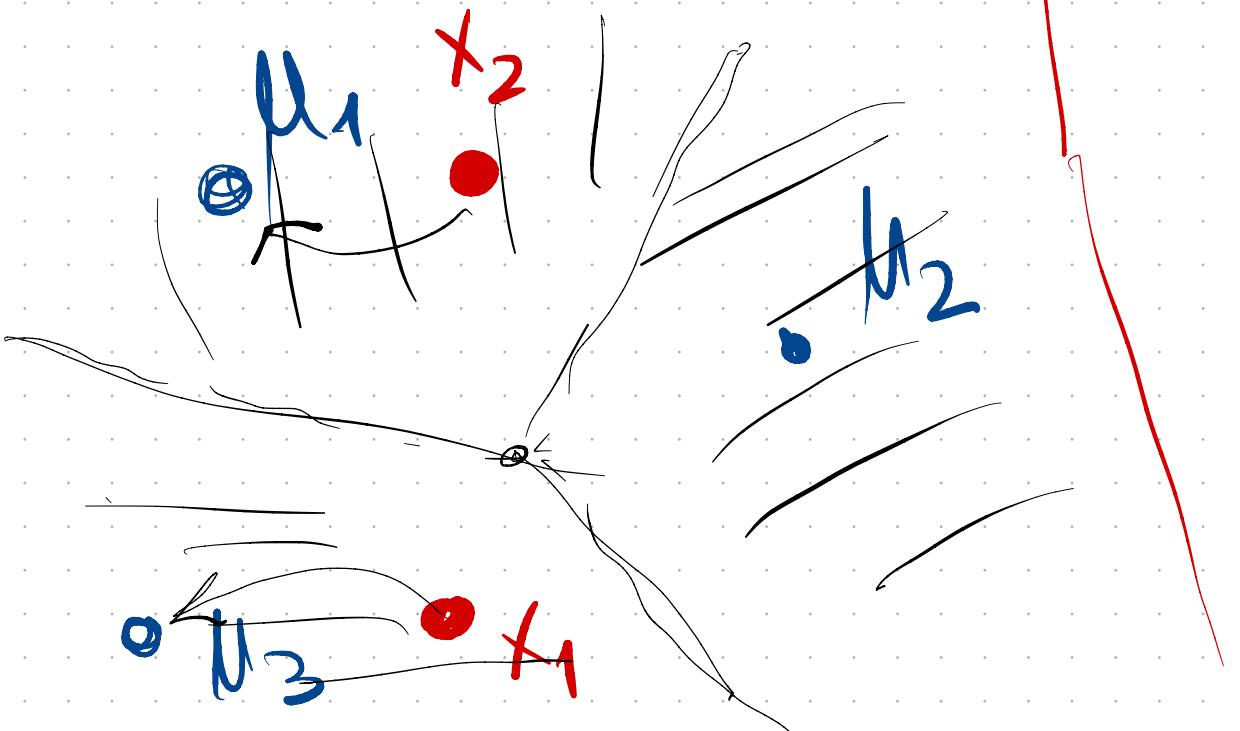
$$= \frac{\sum_{i=1}^N x_i \cdot \pi_{ik}}{\sum_{i=1}^N \pi_{ik}}$$

filter
 only points in
 cluster k

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

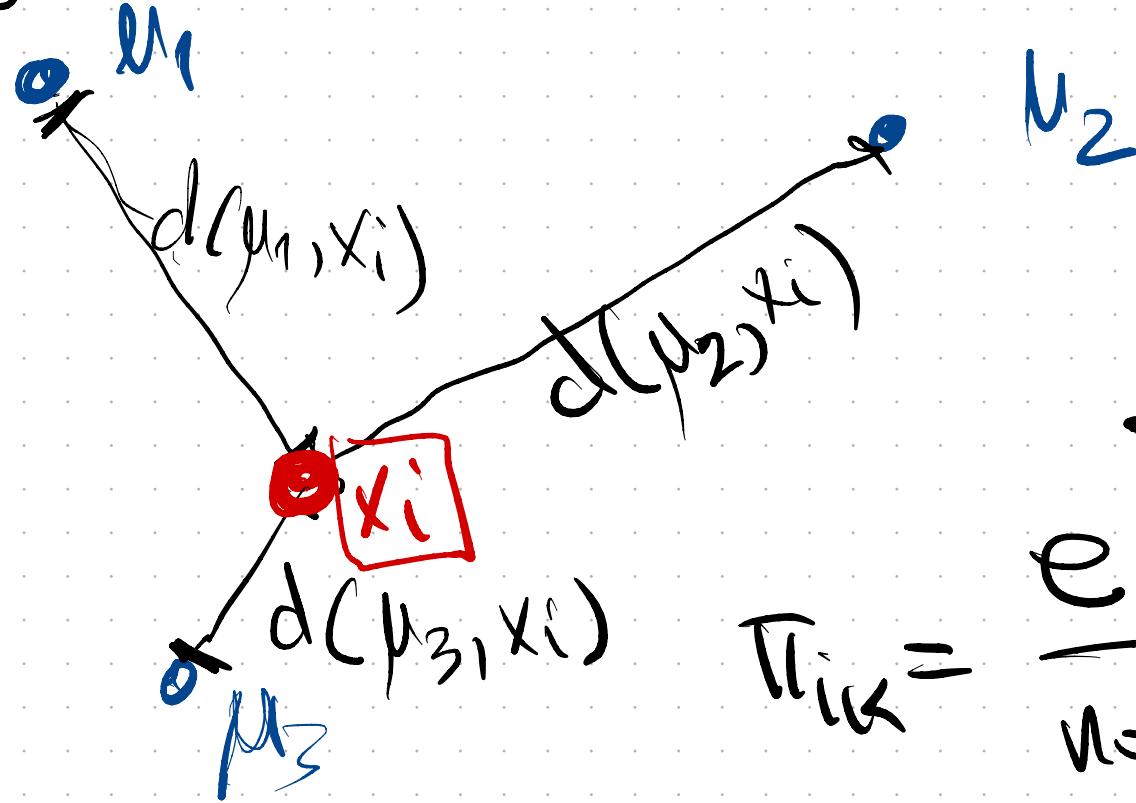
$= N_k = \# \text{ of points}$
 in cluster k

HARD EASY



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

- score dist reg output \Rightarrow probability



M step (same)

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

proportion (weight)

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

$$\pi_{ik} = \frac{e^{-\gamma_p \frac{d(x_i, \mu_k)}{\|x_i - \mu_k\|^2}}}{\text{normalized over all } k}$$

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

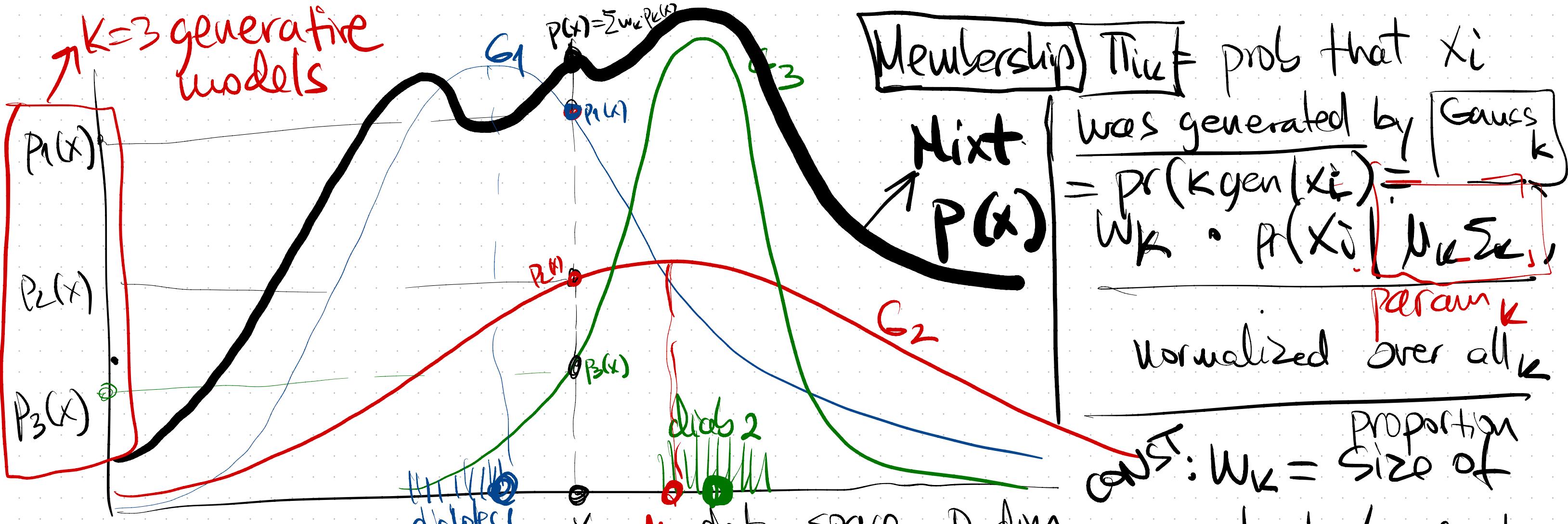
$[\pi_{i1}, \pi_{i2}, \dots, \pi_{iK}]$ distribution over clusters

- parts of x_i distributed $\rightarrow x_i \cdot \pi_{ik}$

• which cluster/group x_i is in

which generator μ_k produced datap. x_i

$\nearrow K=3$ generative models



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

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

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

Membership: Thinf prob that x_i
was generated by Gauss k

$$= \frac{\Pr(k \text{ gen } x_i)}{w_k \cdot P(x_i)} = \frac{w_k \cdot P(x_i | \mu_k, \Sigma_k)}{w_k \cdot P(x_i)}$$

param k

Normalized over all k

const: $w_k = \frac{\text{proportion}}{\text{size of cluster / source } k}$

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

$$= \left[\sum_{i=1}^N \pi_{ik} \right] / N$$

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) given x_i

E step

calculate π_{ik} given

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

probs of x_i generated by comp-k

$$\langle \pi_{ik} \rangle =$$

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

$$= \text{pr}(K | x_i) =$$

$$\frac{\Pr(x_i | k) \cdot \Pr(k)}{\Pr(x_i)}$$

$$\mathcal{N}_k(x_i | \mu_k, \Sigma_k)$$

= normalize over K generators

$$\bullet w_k$$

M step

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

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

- take expectation w.r.t π_{ik}

$$E[LL(x)] \approx LL(x)$$

- take differentials

$$\frac{\partial LL(x)}{\partial \mu_k} \Rightarrow \frac{\partial LL(x)}{\partial \Sigma_k} = 0$$

constraint
 $w_1 + w_2 + w_3 = 1$

Lagrangian

$$\frac{\partial LL(x)}{\partial w_k} \Rightarrow ?$$

$$\text{Likelihood}(\text{data } x) = \prod_{i=1}^n P(x_i) = \prod_{i=1}^n \prod_{k=1}^K [P(x_i | \pi_{ik}) \cdot P(\pi_{ik})]$$

\$\pi_{ik}\$
\$R_k\$
indep of \$i\$
filter
\$1/0\$

$$\text{Log likelihood}(x) = \sum_{i=1}^n \sum_{k=1}^K [\pi_{ik} \log(p(x_i | \pi_{ik})) + \pi_{ik} \log(w_k)]$$

TOO HARD TO OPTIMIZE! $\Rightarrow E[\text{log likelihood}] \text{ w.r.t. } \pi_{ik}$

$$ELL(x) = \sum_{i=1}^n \sum_{k=1}^K \langle \pi_{ik} \log N_k(x_i | \mu_k, \Sigma_k) + \pi_{ik} \log(w_k) \rangle_{\text{prob}}$$

(avg over \$\pi_{ik}\$)

$$N_k(x_i | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \exp \left\{ -\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right\}$$

→ dim gauss

① Optimize for \$\mu_k\$: $\log(N_k(x_i | \mu_k, \Sigma_k)) =$

$$= \log \frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} - \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)$$

$$\frac{\partial \log(\mathcal{W}_K(x_i))}{\partial \mu_k} = \frac{\partial}{\partial \mu_k} \left[\log \frac{1}{(2\pi)^{D/2} (\Sigma_K)^{1/2}} \right] - \frac{\partial}{\partial \mu_k} \left[\frac{1}{2} (x_i - \mu_k) \Sigma_k^{-1} (x_i - \mu_k)^T \right]$$

one point

$$= \frac{1}{2} \frac{\partial}{\partial \mu_k} (x_i - \mu_k) \Sigma_k^{-1} (x_i - \mu_k)$$

$$= (x_i - \mu_k) \Sigma_k^{-1}$$

Vector diff: $\frac{\partial}{\partial x} [x A x^T] = x (A + A^T)$

$$\frac{\partial \text{ELL}(x)}{\partial \mu_k} = \sum_{i=1}^N \langle \pi_{ik} \rangle [(x_i - \mu_k) \Sigma_k^{-1}]$$

$$\Sigma_k^{-1} + (\Sigma_k^{-1})^T = 2\Sigma_k^{-1}$$

all data

want = 0

$$\sum_{i=1}^N \langle \pi_{ik} \rangle x_i = \sum_{i=1}^N \mu_k \cdot \langle \pi_{ik} \rangle$$

$$\Rightarrow \mu_k = \frac{\sum_{i=1}^N \langle \pi_{ik} \rangle x_i}{\sum_{i=1}^N \langle \pi_{ik} \rangle}$$

only k term relevant

$\times \sum_{i=1}^N$
weighted Avg

② Diff wrt $(\Sigma_k)^{-1}$ = COVAR MATRIX OF component k

$$\frac{\partial \log \pi_k(x_i | \mu_k, \Sigma_k)}{\partial \Sigma_k^{-1}} = \frac{\partial}{\partial \Sigma_k^{-1}} \left\{ \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x_i - \mu_k) \Sigma_k^{-1} (x_i - \mu_k)^T \right\}$$

$$= \frac{1}{2} \sum_{k=1}^K - \frac{1}{2} (x_i - \mu_k)^T (x_i - \mu_k)$$

$D \times D$ $1 \times D$ $D \times 1$

All datapoints!

$$\sum_{i=1}^N \langle \pi_{ik} \rangle \left[\frac{1}{2} \Sigma_k - \frac{1}{2} (x_i - \mu_k)^T (x_i - \mu_k) \right] = 0$$

$$\sum_{i=1}^N \langle \pi_{ik} \rangle \Sigma_k = \sum_{i=1}^N \langle \pi_{ik} \rangle (x_i - \mu_k)^T (x_i - \mu_k)$$

$$\Rightarrow \Sigma_k = \frac{1}{\sum_{i=1}^N \langle \pi_{ik} \rangle} \cdot \boxed{(x_i - \mu_k)^T (x_i - \mu_k)}$$

matrix diff

$$\frac{\partial \log M}{\partial M} = (M^{-1})^T$$

$$\frac{\partial}{\partial M}$$

$$\frac{\partial}{\partial M} X M X^T = X^T X$$

EMPIRICAL
COVAR (comp k)

$$\sum_{i=1}^N \boxed{\langle \pi_{ik} \rangle} \text{ weighted avg}$$

③ diff w.r.t w_k = mixture coef constraint $w_1+w_2+w_3=1$

Lagrangian $\text{ELL}(x_i) - \lambda \cdot [w_1+w_2+w_3-1]$

$$\frac{\partial \text{Lag}}{\partial w_k} = \frac{\partial}{\partial w_k} \left[\langle \pi_{ik} \rangle \cdot \log(w_k) - \lambda [w_1+w_2+w_3-1] \right]$$

$$= \langle \pi_{ik} \rangle \cdot \frac{1}{w_k} - \lambda$$

All datapoints: $\frac{\partial \text{Lag}}{\partial w_k} = \sum_{i=1}^N \langle \pi_{ik} \rangle \cdot \frac{1}{w_k} - \lambda = 0$ want

$$\sum_{i=1}^N \langle \pi_{ik} \rangle = \lambda w_k \Rightarrow w_k =$$

$$\frac{\sum_{i=1}^N \langle \pi_{ik} \rangle}{\lambda} = \frac{\sum_{i=1}^N \langle \pi_{ik} \rangle}{N}$$

constraint $1=w_1+w_2+w_3 \Rightarrow$
 (or equiring $\frac{\partial}{\partial \lambda}$)

$$\boxed{\sum_{i=1}^N \sum_{k=1}^K \langle \pi_{ik} \rangle / \lambda = 1}$$

$$\Rightarrow \boxed{\lambda = N}$$