

GMM-EM

4/29/26

(1)

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

k = cluster index

π_k = selection prob \rightarrow mixing coefficients

μ_k, Σ_k = single gaussian params.

θ = some params = $\{\pi_k, \mu_k, \Sigma_k\}$

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

$$Z(x|\theta) = \prod_{n=1}^N p(x_n|\theta)$$

$$\mathcal{L} = \log Z(x|\theta) = \log \prod_{n=1}^N p(x_n|\theta)$$

$$\Rightarrow \sum_{n=1}^N \log p(x_n|\theta)$$

$$\Rightarrow \sum_{n=1}^N \log \left[\sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k) \right]$$

~~instead of weights, assume we know which gaussian a datapoint comes from.~~

Log-likelihood of single datapoint.

$$\ln p(x|\theta) = \sum_z q(z) \ln p(x|\theta)$$

Introduce additional variable for fun-sies..

Now rewrite $p(x|\theta)$ to include z .

$$p(x|\theta) = \frac{p(x, \theta)}{p(\theta)} = \frac{\sum_z p(x, \theta, z)}{\sum_z p(\theta, z)} = \frac{\sum_z p(x, z|\theta) p(\theta)}{\sum_z p(\theta, z)} \quad ??$$

see next page

(2)

try working backwards,

$$\text{known: } q(z|\theta) = \frac{p(x, z|\theta)}{p(z|x, \theta)} = \frac{p(x, z, \theta) / p(\theta)}{p(x, z, \theta) / p(x, \theta)} = \frac{1/p(\theta)}{1/p(x, \theta)}$$

$$\Rightarrow \frac{p(x, \theta)}{p(\theta)} = p(x|\theta)$$

Inverse,

$$p(x|\theta) = \frac{p(x, \theta)}{p(\theta)} = \frac{p(x, z, \theta)}{p(x, z, \theta)} \cdot \frac{p(x, \theta)}{p(\theta)}$$

$$\Rightarrow \frac{p(x, z, \theta)}{p(\theta)} \cdot \left[\frac{p(x, \theta)}{p(x, z, \theta)} \right]$$

$$\Rightarrow p(x, z|\theta) \cdot \left[\frac{p(x, z, \theta)}{p(x, \theta)} \right]^{\uparrow}$$

$$\Rightarrow p(x, z|\theta) \cdot \left[p(z|x, \theta) \right]^{\uparrow}$$

$$\Rightarrow \frac{p(x, z|\theta)}{p(z|x, \theta)}$$

$$\Rightarrow \ln p(x|\theta) = \sum_z q(z) \ln p(x|\theta) = \sum_z q(z) \ln \left[\frac{p(x, z|\theta)}{p(z|x, \theta)} \right]$$

$$\Rightarrow \sum_z q(z) \ln \left[\frac{p(x, z|\theta)}{p(z|x, \theta)} \cdot \frac{q(z)}{q(z)} \right] \rightarrow \text{for Jensen's.}$$

$$\Rightarrow \sum_z q(z) \ln \left[\frac{p(x, z|\theta)}{q(z)} \cdot \frac{q(z)}{p(z|x, \theta)} \right]$$

$$\Rightarrow \sum_z q(z) \left[\ln \left[\frac{p(x, z|\theta)}{q(z)} \right] + \ln \left[\frac{q(z)}{p(z|x, \theta)} \right] \right]$$

$$\Rightarrow \sum_z q(z) \ln \left[\frac{p(x, z|\theta)}{q(z)} \right] + \sum_z q(z) \ln \left[\frac{q(z)}{p(z|x, \theta)} \right]$$

$$\ln p(x|\theta) = \underbrace{\sum_z q(z) \ln \left[\frac{p(x, z|\theta)}{q(z)} \right]}_{\substack{\mathcal{L}(q, \theta) \\ \text{ELBO}}} + \underbrace{\sum_z q(z) \ln \left[\frac{q(z)}{p(z|x, \theta)} \right]}_{\text{KL}(q||p)}$$

$\text{KL}(q||p) \geq 0$ due to Jensen's Inequality:

$$\text{KL}(q||p) = \sum_n q(x) \ln \left[\frac{q(x)}{p(x)} \right] = -\sum_n q(x) \ln \left[\frac{p(x)}{q(x)} \right] \Rightarrow -\mathbb{E}_q \left[\ln \left[\frac{p(x)}{q(x)} \right] \right]$$

Jensen's Inequality for convex functions:

Expectation $\mathbb{E}(f(u)) \geq f(\mathbb{E}(u))$

$f = -\ln$ only as $+\ln$ is not convex

$$\mathbb{E}_q [\ln(u)] \geq -\ln[\mathbb{E}(u)]$$

~~$$\mathbb{E}_q [\ln(u)] \geq -\ln[\mathbb{E}(u)] \Rightarrow \mathbb{E}_q \left[\ln \left[\frac{p(x)}{q(x)} \right] \right] \geq -\ln \left[\mathbb{E}_q \left[\frac{p(x)}{q(x)} \right] \right]$$~~

$$-\sum_n q(x) \ln \left[\frac{p(x)}{q(x)} \right] \geq -\ln \left[\sum_n q(x) \frac{p(x)}{q(x)} \right] \geq -\ln \left[\sum_n p(x) \right]$$

$$\boxed{-\sum_n q(x) \ln \left[\frac{p(x)}{q(x)} \right] \geq 0}$$

$$p(x, z|\theta) = \frac{p(x, z, \theta)}{p(\theta)}$$
~~$$p(x, \theta|z) p(z)$$~~
~~$$p(x, \theta, z) p(z)$$~~

$$\Rightarrow \frac{p(x|z, \theta) p(z, \theta)}{p(\theta)}$$

$$\Rightarrow p(x|z, \theta) \frac{p(z, \theta)}{p(\theta)}$$

$$\Rightarrow p(x|z, \theta) p(z|\theta)$$

$$p(x|\theta) = \sum_{k=1}^K p(x, z_k=1|\theta)$$

cluster selector.

~~$$p(x|\theta)$$~~

$$p(x, z|\theta) = p(z|\theta) p(x|z, \theta)$$

prior over latent cluster.

(4)

$$\ln p(x|\theta) = \log \left[\sum_{k=1}^K \pi_k e^{\sqrt{(x_k | \mu_k, \Sigma_k)}} \right] = \log \sum_{k=1}^K p(x, z_k | \theta)$$

~~Instead if we do marginalize z_k (or z_k)~~

$$\Rightarrow \log \sum_{k=1}^K p(z_k | \theta) p(x | z_k, \theta) \Rightarrow \log \sum_{k=1}^K p(x, z_k | \theta)$$

Instead of marginalizing over z_k , we can introduce an alternative form for $p(x|\theta) = \frac{p(x, z | \theta)}{p(z | x, \theta)}$ which leads

us to the ELBO + KL decomposition.

$$\Rightarrow \ln p(x|\theta) = \sum_z q(z) \ln \left[\frac{p(x, z | \theta)}{q(z)} \right] + \sum_z q(z) \ln \left[\frac{q(z)}{p(z | x, \theta)} \right]$$

~~$$\sum_z q(z) \ln \left[\frac{p(x, z | \theta)}{q(z)} \right] = \sum_z q(z) \ln \left[\prod_k \pi_k e^{\sqrt{(x_k | \mu_k, \Sigma_k)}} \right]$$~~

$$p(x, z | \theta) = \prod_{k=1}^K \left[\pi_k e^{\sqrt{(x_k | \mu_k, \Sigma_k)}} \right]^{z_k} \text{ for } z_k = \delta(\text{cluster} = k)$$

EM Intuition : $\ln p(x|\theta) = \underbrace{2(q|\theta)}_{\text{constant for some } \theta} + \underbrace{KL_{q||p}}_{\geq 0}$

If we hold θ fixed : $\Delta \ln p(x|\theta) = \Delta 2(q|\theta) + \Delta KL_{q||p}$

$$\Delta 2(q|\theta) = - \Delta KL_{q||p}$$

If we minimize $KL_{q||p}$ for θ fixed, the ELBO will go up.

$$KL_{q||p} = 0 \text{ for } q = p.$$

$$\Rightarrow q(z) = p(z | x, \theta \text{ fixed})$$

w/ $KL_{q||p} = 0$, we have $\ln p(x|\theta) = 2(q, \theta)$
 and because $KL_{q||p} = 0$ and can't go below 0,
 any changes we make to $2(q, \theta)$ cannot reduce
 $\ln p(x|\theta)$.

E-step

Start w/ find q such that $KL_{q||p} = 0$

$\Rightarrow q(z) = p(z|x, \theta) \Rightarrow p(z|x)$ as $\theta \rightarrow$ fixed

$p(z|x) p(x) = p(x|z) p(z)$

1 datapoint $\left[\Rightarrow p(z|x) = \frac{p(x|z) p(z)}{p(x)} = \frac{\sum_k \pi_k e^{V(x|\mu_k, \Sigma_k)} \pi_k}{\sum_k \pi_k e^{V(x|\mu_k, \Sigma_k)}} \right]$
 Marginalized $p(x, z)$

all points $\Rightarrow q(z) = \prod_{n=1}^N q(z_n)$

$q(z) \Rightarrow \prod_{n=1}^N \prod_{k=1}^K [\gamma_{nk}]^{z_{nk}}$

N datapoint $= \prod_{n=1}^N \frac{\sum_k \pi_k e^{V(x_n|\mu_k, \Sigma_k)} \pi_k}{\sum_k \pi_k e^{V(x_n|\mu_k, \Sigma_k)}}$

over all z
 $p(z_n|x) = \prod_{k=1}^K [\gamma_{nk}]^{z_{nk}} = q(z_n)$

i.e. responsibility for all clusters for a datapoint

(6)

M-step

Insert $q(z)$ [such that $q(z) = p(z|\alpha, \theta)$] into complete data log-likelihood wh/ makes $K_{\theta} q_{\theta} p = 0$

$$\ln p(\alpha|\theta) = \sum_z q(z) \ln \left[\frac{p(\alpha, z|\theta)}{q(z)} \right]$$

$$\Rightarrow \sum_z q(z) [\ln p(\alpha, z|\theta) - \ln q(z)]$$

$$\Rightarrow \underbrace{\sum_z q(z) \ln p(\alpha, z|\theta)}_{\mathbb{E}_q[\ln p(\alpha, z|\theta)]} - \underbrace{\sum_z q(z) \ln q(z)}_{\text{Entropy}}$$

determined using $\theta^{\text{old}} = \text{const}$
 $\therefore \frac{\partial}{\partial \theta} = 0$

~~$\ln p(\alpha|\theta) = \sum_z q(z) \ln p(\alpha, z|\theta)$~~

$$\begin{aligned} & \ln \prod_n p(\alpha_n, z_n|\theta) \\ &= \sum_n \ln p(\alpha_n, z_n|\theta) \\ &= \sum_n \ln \prod_k [\pi_k \mathcal{N}(\alpha_n | \mu_k, \Sigma_k)]^{z_{nk}} \\ &= \sum_n \sum_k \ln [\pi_k \mathcal{N}(\alpha_n | \mu_k, \Sigma_k)]^{z_{nk}} \end{aligned}$$

$$\Rightarrow \mathbb{E}_q \left[\sum_n \sum_k \ln [\pi_k \mathcal{N}(\alpha_n | \mu_k, \Sigma_k)]^{z_{nk}} \right]$$

$$\Rightarrow \mathbb{E}_q \left[\sum_n \sum_k z_{nk} \ln [\pi_k \mathcal{N}(\alpha_n | \mu_k, \Sigma_k)] \right] \quad \mathbb{E}(A+B) = \mathbb{E}(A) + \mathbb{E}(B)$$

$$\Rightarrow \sum_n \sum_k \underbrace{\mathbb{E}_q \left[z_{nk} \ln [\pi_k \mathcal{N}(\alpha_n | \mu_k, \Sigma_k)] \right]}_{\text{Constant w.r.t } q}$$

$$\Rightarrow \sum_n \sum_k \mathbb{E}_q(z_{nk}) \ln [\pi_k \mathcal{N}(\alpha_n | \mu_k, \Sigma_k)]$$

$$\Rightarrow \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}_q(z_{nk}) \ln [\pi_k \mathcal{N}(x_k | \mu_k, \Sigma_k)]$$

$$\begin{aligned} \Rightarrow \mathbb{E}_q(z_{nk}) &= \sum_{z_n} q(z_n) z_{nk} \quad \leftarrow \text{selector variable} \\ &= [1 \cdot q(z_{nk}=1)] + [0 \cdot q(z_{nk}=0)] \\ &= q(z_{nk}=1) = \gamma_{nk} \end{aligned}$$

$$\Rightarrow \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \ln [\pi_k \mathcal{N}(x_k | \mu_k, \Sigma_k)] = \ln p(x | \theta)$$

Derive parameter updates:

of a specific cluster so don't have to sum over k.

$$\frac{\partial}{\partial \mu_k} \ln p(x | \theta) = \frac{\partial}{\partial \mu_k} \sum_{n=1}^N \gamma_{nk} \left[\ln \pi_k + \ln \mathcal{N}(x_k | \mu_k, \Sigma_k) \right] = Q$$

$$\Rightarrow \frac{\partial}{\partial \mu_k} \sum_{n=1}^N \gamma_{nk} \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right] \quad \left| \frac{\partial}{\partial \mu_k} \ln \left[\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right) \right] \right.$$

$$\Rightarrow \sum_{n=1}^N \gamma_{nk} \cdot \frac{\partial}{\partial \mu_k} \left[\underbrace{(x - \mu)^T}_{u} \Sigma^{-1} \underbrace{(x - \mu)}_v \right] \quad \Rightarrow \ln \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)$$

$$\frac{\partial u}{\partial \mu_k} = -1, \quad \frac{\partial v}{\partial \mu_k} = -\Sigma^{-1}$$

$$\frac{\partial}{\partial \mu_k} (u \cdot v) = u'v + uv' \quad \Rightarrow -\Sigma^{-1} (x - \mu) - (x - \mu)^T \Sigma^{-1}$$

since Σ^{-1} is symmetric $[a \Sigma^{-1}]^T = \Sigma^{-1} a^T$

$$\Rightarrow -2 \Sigma^{-1} (x - \mu)$$

$$\frac{\partial Q}{\partial \mu_k} = \frac{1}{2} \sum_{n=1}^N \gamma_{nk} [-2 \Sigma^{-1} (x - \mu)] = \sum_{n=1}^N \gamma_{nk} \Sigma_k^{-1} (x_k - \mu_k)$$

8

find max

$$\frac{\partial Q}{\partial \mu_k} = 0 = \sum_n \gamma_{nk} \Sigma_k^{-1} (x_n - \mu_k) \Rightarrow \sum_n \gamma_{nk} \Sigma_k^{-1} x_n = \sum_n \gamma_{nk} \Sigma_k^{-1} \mu_k$$

$$\Rightarrow \frac{\sum_n \gamma_{nk} \Sigma_k^{-1} x_n}{\sum_n \gamma_{nk} \Sigma_k^{-1}} = \mu_k = \frac{\sum_n \sum_k \gamma_{nk} x_n}{\sum_n \sum_k \gamma_{nk}}$$

$$= \mu_k \sum_n \gamma_{nk} \Sigma_k^{-1}$$

$$\mu_k = \frac{\sum_n \gamma_{nk} x_n}{\sum_n \gamma_{nk}}$$

$$\frac{\partial Q}{\partial \pi_k} = \sum_n \gamma_{nk} \ln \pi_k$$

But we want to optimize subject to constraint $\sum_k \pi_k = 1$

Constrained optimization:

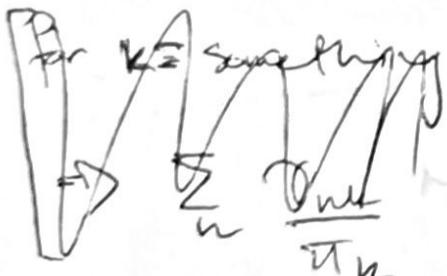
$$\Rightarrow J = \sum_n \sum_k \gamma_{nk} \ln \pi_k + \lambda (\sum_k \pi_k - 1)$$

$$\frac{\partial J}{\partial \pi_k} = \sum_n \sum_k \gamma_{nk} \cdot \frac{1}{\pi_k} + \frac{\partial}{\partial \pi_k} \left[\lambda \sum_k \pi_k + \lambda \sum_k (-1) \right]$$

$$\Rightarrow \sum_n \sum_k \frac{\gamma_{nk}}{\pi_k} + \sum_k \lambda = 0$$

$$\frac{\partial}{\partial \pi_k} = \lambda_k \quad \frac{\partial}{\partial \pi_k} = 0$$

$$\text{or } \frac{\partial}{\partial \pi_k} = \sum_n \lambda$$



$$\Rightarrow \sum_k \left[\sum_n \frac{\gamma_{nk}}{\pi_k} + \lambda \right] \Rightarrow \sum_k \left[\sum_n \gamma_{nk} + \lambda \pi_k \right]$$

$$\text{for the case } k \Rightarrow \sum_n \gamma_{nk} + \lambda \pi_k = 0 \Rightarrow \sum_n \gamma_{nk} = -\lambda \pi_k$$

$$\Rightarrow \sum_n \sum_k \frac{\gamma_{nk}}{\pi_k} = \sum_n \sum_k \frac{\gamma_{nk}}{\pi_k}$$

$$\lambda = \frac{\sum_n \gamma_{nk}}{\pi_k}$$

$$\frac{\partial \bar{J}}{\partial \pi_k} = \sum_n \sum_k \frac{\gamma_{nk}}{\pi_k} + \sum_k \lambda = 0$$

$$\left[\sum_k \left[\frac{1}{\pi_k} \sum_n \gamma_{nk} \right] + \sum_k \lambda = 0 \Rightarrow \sum_k \left[\frac{1}{\pi_k} \sum_n \gamma_{nk} + \lambda \right] = 0 \right]$$

∴ can lose the \sum_k b/c only the k -th term will have a non-zero derivative.

$$\Rightarrow \sum_n \frac{\gamma_{nk}}{\pi_k} + \lambda = 0 \Rightarrow \lambda = - \sum_n \frac{\gamma_{nk}}{\pi_k}$$

$$\Rightarrow \pi_k \lambda = - \sum_n \gamma_{nk} \Rightarrow \lambda \sum_k \pi_k = - \sum_n \sum_k \gamma_{nk}$$

All responsibility for any point $\Sigma=1$

$$\Rightarrow \lambda = - \sum_k 1 \Rightarrow \boxed{\lambda = -N}$$

$$\sum_n \frac{\gamma_{nk}}{\pi_k} + (-N) = 0 \Rightarrow \sum_n \frac{\gamma_{nk}}{\pi_k} = N \Rightarrow \boxed{\frac{1}{N} \sum_n \gamma_{nk} = \pi_k}$$

\sum_k

$$\frac{\partial}{\partial \Sigma_k} \ln p(x|\theta) = \frac{\partial}{\partial \Sigma_k} \sum_n \gamma_{nk} \left[\ln \pi_k + \ln \mathcal{N}(x_n | \mu, \Sigma_k) \right]$$

$$\frac{\partial}{\partial \Sigma_k} = 0$$

$$\Rightarrow \sum_n \gamma_{nk} \frac{\partial}{\partial \Sigma_k} \left[\ln \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{d/2}} - \frac{1}{2} (x-\mu)^T \Sigma_k^{-1} (x-\mu) \right]$$

$$\Rightarrow \frac{\partial}{\partial \Sigma_k} \ln [(2\pi)^{d/2} |\Sigma_k|^{d/2}]^{-1} \Rightarrow \frac{\partial}{\partial \Sigma_k} \ln [(2\pi)^{d/2} |\Sigma_k|]^{-1/2} = \frac{\partial}{\partial \Sigma_k} -\frac{1}{2} \ln [(2\pi)^d |\Sigma_k|]$$

$$\Rightarrow -\frac{1}{2} \frac{\partial}{\partial \Sigma_k} [\ln (2\pi)^d + \ln |\Sigma_k|] \Rightarrow -\frac{1}{2} \Sigma^{-1}$$

via magic

(10)

$$\frac{\partial}{\partial \Sigma_k} \underbrace{(x-\mu)^T}_{u} \underbrace{\Sigma_k^{-1}}_v (x-\mu) \Rightarrow \frac{\partial u}{\partial \Sigma_k} (x-\mu)^T \Sigma_k^{-1} \Rightarrow (x-\mu)^T \frac{\partial}{\partial \Sigma_k} \Sigma^{-1}$$

$$\frac{\partial}{\partial \Sigma_k} u v = u' v + u v' \quad \frac{\partial v}{\partial \Sigma_k} = 0$$

$$\Rightarrow (x-\mu)^T \left[\Sigma^{-1} \frac{d\Sigma \Sigma^{-1}}{d\Sigma} \right]$$

$$\Rightarrow (x-\mu)^T \left[-\Sigma^{-1} \frac{d\Sigma \Sigma^{-1}}{d\Sigma} \right] (x-\mu) \stackrel{\text{by magic}}{\Rightarrow} -\Sigma^{-1} (x-\mu)(x-\mu)^T \Sigma^{-1}$$

$$\frac{\partial}{\partial \Sigma_k} \ln p(x|\theta) = \sum_{n=1}^n \gamma_{nk} \left[\frac{1}{2} \Sigma_k^{-1} + \frac{1}{2} \Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} \right]$$

$$\Rightarrow \frac{1}{2} \sum_{n=1}^n \gamma_{nk} \left[\Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} - \Sigma_k^{-1} \right]$$

$$= \frac{\Sigma_k^{-1}}{2} \sum_{n=1}^n \gamma_{nk} \left[(x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} - \mathbb{1} \right]$$

$$= \frac{\Sigma_k^{-1}}{2} \left[\sum_{n=1}^n \gamma_{nk} \left[(x_n - \mu_k)(x_n - \mu_k)^T - \Sigma_k \right] \right] \Sigma_k^{-1} = 0$$

$$\Rightarrow \sum_{n=1}^n \gamma_{nk} \left[(x_n - \mu_k)(x_n - \mu_k)^T - \Sigma_k \right] = 0$$

$$\Rightarrow \sum_{n=1}^n \gamma_{nk} \left[(x_n - \mu_k)(x_n - \mu_k)^T \right] - \sum_{n=1}^n \gamma_{nk} \Sigma_k = 0$$

$$\sum_{n=1}^n \gamma_{nk} \left[(x_n - \mu_k)(x_n - \mu_k)^T \right] = \sum_{n=1}^n \gamma_{nk} \Sigma_k = \Sigma_k \sum_{n=1}^n \gamma_{nk}$$

$$\boxed{\Sigma_k = \frac{\sum_{n=1}^n \gamma_{nk} \left[(x_n - \mu_k)(x_n - \mu_k)^T \right]}{\sum_{n=1}^n \gamma_{nk}}}$$