

Gaussian Mixture Models

jhsh
jhshi@vt.edu

April 8, 2021

Introduction: Gaussian as Mixtures

Model: Data $x \in \mathbf{R}^M$ disbursed *normally* with K centroids

$$p(x_i | (\mu, \sigma)_k) = \mathcal{N}(x_i, (\mu, \sigma)_k) \equiv g_{ik} \quad x_i = \text{Observable Data Point } i$$
$$(\mu, \sigma)_k = \text{mean, covariance } k$$

Hidden DOF $(z_j) \equiv$ "one hot" at some $k := 1(k)$ **k-th centroid**

$$p(z_j | \pi) = \prod_{s \in [1, \mathcal{K}]} \pi_s^{(z_j)^s} \quad \text{"simplex:"} \quad \sum_{s \in [1, \mathcal{K}]} \pi_s = 1$$

$$p(x_i | z_i(\mu, \sigma)) = \prod_{t \in [1, \mathcal{K}]} g_{it}^{(z_i)^t} \quad p(x | z(\mu, \sigma)) = \prod_{i \in [1, N]} p(x_i | z_i(\mu, \sigma))$$

Introduction: z_j unobserved vs. observed

$\mathcal{L}_N(\theta | x, z) \equiv$ "log likelihood" : $\arg \max_{\theta} p(x, z | \theta) \simeq \arg \max_{\theta} \mathcal{L}_N(\theta | x, z)$

- if z_j remained *hidden*:

$$\mathcal{L}(\theta, x, z) = \sum_{i \in [1, N]} \log \sum_{z_i} p(x_i | z_i \theta) p(z_i | \theta) = \sum_{i \in [1, N]} \log \sum_{z_i} \prod_{s \in [1, \mathcal{K}]} (\pi_s g_{is})^{(z_i)^s}$$

\sum_{z_i} **Hard to Compute...**

- assume *complete* observability $D = (x, z)$: $z_j^k := \mathbf{1}(k)$ k -th cluster

$$\mathcal{L}_N(\theta | x, z) = \sum_{i \in [1, N]} \sum_{j \in [1, \mathcal{K}]} (z_i)^j \cdot \log(\pi_j g_{ij})$$

Q-Learning(?)

convenient to define $Q(\theta^{(t)}) \equiv \langle \mathcal{L}(\theta, x, z) \rangle_{p \sim z | x \theta^{(t)}}$ **iteration t**

$$\begin{aligned} Q(\theta^{(t+1)}) &\leftarrow \left\langle \sum_{i \in [1, N]} \sum_{j \in [1, \mathcal{K}]} (z_i)^j \log(\pi_j g_{ij})^{(t)} \right\rangle_{p \sim z | x \theta^{(t)}} \\ &= \sum_{i \in [1, N]} \sum_{j \in [1, \mathcal{K}]} r_{ij}^{(t)} \log(\pi_j g_{ij})^{(t)} \end{aligned}$$

$$r_{ij} = \langle (z_i)^j \rangle_{p \sim z | x \theta^{(t)}} = p((z_i)^j | \theta, x_i) = (\pi_k g_{ik}) (\sum_{s \in [1, \mathcal{K}]} (\pi_s g_{is}))^{(-)}$$

"responsibilities"

Optimization:

- update (r, Q) :

$$r_{ik}^{(t+1)} \leftarrow (\pi_k g_{ik}) \left(\sum_{s \in [1, \mathcal{K}]} (\pi_s g_{is})^{(t)} \right)^{(-)}$$

- update parameters θ :

$$\pi_k^{t+1} \leftarrow N^{(-)} \left(\sum_{i \in [1, N]} r_{ik}^t \right)$$

$$\mu_k^{t+1} \leftarrow \left(\sum_{i \in [1, N]} r_{ik}^t \right)^{(-)} \left(\sum_{i \in [1, N]} r_{ik}^t x_i \right)$$

$$\sigma_k^{t+1} \leftarrow \frac{\sum_i r_{ik} \|x_i - \mu_k\|^2}{2 \sum_i r_{ik}}$$

eg: $0 \equiv \frac{\partial}{\partial \pi_k} \left(Q(\theta^{(t)}) + \lambda (\sum_{s \in [1, \mathcal{K}]} \pi_s - 1) \right) \rightarrow \lambda = -N$

*** derivations ***

Blank Page

review: EM

- Let x be observed data and θ be the underlying parameters, and z be the latent variables.
- "incomplete" log-likelihood: $\mathcal{L}_N(x, \theta) = \log p(x|\theta)$
- "complete" log-likelihood: $\mathcal{L}_N(x, \theta)^{\text{comp}} = \log p(x, z|\theta)$
- But $\mathcal{L}_N(x, \theta)^{\text{comp}}$ is not physically observable, since z is still hidden from view. One proposal, i.e EM, is to trace off z along with its posterior.
- EM algorithm:

- At present t , calculate the Q -function:

$$Q(\theta | \theta') \equiv \langle \log p(x, z, \theta) | x\theta' \rangle \quad \text{where} \quad \langle \bullet | x\theta' \rangle \sim \sum_z p(z | x\theta')(\bullet) \quad \text{"E Step"}$$

- maximize Q with respect to θ :

$$\theta_t \leftarrow \arg \max_{\theta_t} Q(\theta_{t+1} | \theta_t) \quad \text{"M Step"}$$

- "Likelihood never decreases in EM." $\mathcal{L}_N(x, \theta_t) \leq \mathcal{L}_N(x, \theta_{t+1})$
i.e non-negative rest term: $\Delta(x, \theta, q) \equiv \mathcal{L}_N(x, \theta) - F(\theta, x, q)$
 - $\Delta(x, \theta, q) = D_{\text{KL}}(q \| w)$

continued

- Rate of Convergence: assume $\lim_{t \gg 0} (\theta_t, \theta_{t+1}) \equiv \theta_{(*)}$ and vanishing $(\partial_\theta Q(\theta, \theta'))_{\theta_{(*)}} \sim 0 \sim (\partial_\theta Q(\theta_t, \theta_{t+1}))$

$$R \equiv \lim_{t \gg 0} (\theta_{t+1} - \theta_{(*)})(\theta_t - \theta_{(*)})^{(-)} = \mathcal{J}_{\text{com}}(\theta_{(*)} | x) \mathcal{J}_{\text{cond}}^{(-)}(\theta_{(*)} | x)$$

- Note: $\mathcal{J}_{\text{com}} \equiv (-) \langle \partial_{\theta_*}^2 \log p(z, x, \theta_*) | x, \theta_* \rangle$ complete Fisher metric
 $\mathcal{J}_{\text{cond}} \equiv (-) \langle \partial_{\theta_*}^2 \log p(z | x, \theta_*) | x, \theta_* \rangle$ observed Fisher metric

Derivation

$$R \equiv \lim_{t \gg 0} (\theta_{t+1} - \theta_{(*)})(\theta_t - \theta_{(*)})^{(-)} = \mathcal{J}_{\text{com}}(\theta_{(*)} | x) \mathcal{J}_{\text{cond}}^{(-)}(\theta_{(*)} | x)?$$

- Taylor expand $\lim_{(\theta_1 \theta_2) = (\theta_t \theta_{t+1})} \partial_{\theta_2} Q(\theta_2 | \theta_1)$ around convergent $\theta_{(*)}$:

$$\begin{aligned} \lim_{(\theta_1 \theta_2) = (\theta_t \theta_{t+1})} \partial_{\theta_2} Q(\theta_2 | \theta_1) &\approx \lim_{(\theta_1 \theta_2) = (\theta_t \theta_{t+1})} \left[\lim_{(\theta_1 \theta_2) = (\theta_*)} \partial_{\theta_2} Q(\theta_2 | \theta_1) \right. \\ &\quad + (\theta_1 - \theta_*) \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) + (\theta_2 - \theta_*) \partial_{\theta_2}^2 Q(\theta_2 | \theta_1) \\ &\quad \left. + \mathcal{O}(\sim \theta^2) \right] \end{aligned}$$

- Employ " $\partial_{\theta_2} Q(\theta_2 | \theta_1) \big|_{(\theta_t \theta_{t+1})} = 0 = \partial_{\theta_2} Q(\theta_2 | \theta_1) \big|_{(\theta_* \theta_*)}$ "?

$$\lim_{(\theta_1 \theta_2) = (\theta_t \theta_{t+1})} [(\theta_1 - \theta_*) \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) + (\theta_2 - \theta_*) \partial_{\theta_2}^2 Q(\theta_2 | \theta_1)] = 0$$

$$R = \lim_{t \gg 0} \frac{\theta_{t+1} - \theta_*}{\theta_t - \theta_*} = (-) \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \frac{\partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1)}{\partial_{\theta_2}^2 Q(\theta_2 | \theta_1)}$$

- $\lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) \stackrel{?}{=} (-) \cdot \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1}^2 D(\theta_2 | \theta_1)$

$$\begin{aligned} \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) &= \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1} \partial_{\theta_2} \left(\sum_z p(z | x, \theta_1) \log p(x, z, \theta_2) \right) \\ &= \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1} \partial_{\theta_2} \left(\sum_z p(z | x, \theta_1) \log p(z | x, \theta_2) \right) \end{aligned}$$

- continued:

$$\begin{aligned}
\lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) &= \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \sum_z (\partial_{\theta_1} p(z | x, \theta_1)) \cdot \partial_{\theta_2} \log p(z | x, \theta_2) \\
&= \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \sum_z (\partial_{\theta_*} p(z | x, \theta_*)) \cdot \partial_{\theta_*} \log p(z | x, \theta_*) \\
&= \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \sum_z \partial_{\theta_*} \left[(p(z | x, \theta_*)) \cdot \partial_{\theta_*} \log p(z | x, \theta_*) \right] \\
&\quad + (-) \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \sum_z (p(z | x, \theta_*)) \cdot \partial_{\theta_*}^2 \log p(z | x, \theta_*)
\end{aligned}$$

- The first term vanish under the rule of total derivative.

$$\begin{aligned}
\sum_z \partial_{\theta_*} \left[(p(z | x, \theta_*)) \cdot \partial_{\theta_*} \log p(z | x, \theta_*) \right] &= \sum_z \partial_{\theta_*} \left[p(z | x, \theta_*) \cdot \frac{\partial_{\theta_*} p(z | x, \theta_*)}{p(z | x, \theta_*)} \right] \\
&= \sum_z \partial_{\theta_*} \left[\partial_{\theta_*} p(z | x, \theta_*) \right] = \partial_{\theta_*}^2 \sum_z \left[p(z | x, \theta_*) \right] = \partial_{\theta_*}^2 [1] = 0
\end{aligned}$$

- WLOG: $\lim_{\theta_1 \theta_2 \rightarrow \theta_*} (\bullet) \sim \lim_{t \gg 0} \lim_{\theta_1 \theta_2 \rightarrow \theta_t \theta_{t+1}} (\bullet)$

$$\begin{aligned}
\lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) &= (-) \lim_{(\theta_1 \theta_2) \rightarrow (\theta_* \theta_*)} \sum_z p(z | x, \theta_*) \cdot \partial_{\theta_*}^2 \log p(z | x, \theta_*) \\
&= (-) \lim_{t \gg 0} \sum_z p(z | x, \theta_t) \cdot \partial_{\theta_{t+1}}^2 \log p(z | x, \theta_{t+1}) \\
&= (-) \lim_{t \gg 0} \partial_{\theta_{t+1}}^2 \langle \log p(z | x, \theta_{t+1}) | x, \theta_t \rangle
\end{aligned}$$

- Substitute definitions for various Fisher metrics:

$$\begin{aligned}\lim_{(\theta_1, \theta_2) \rightarrow (\theta_*, \theta_*)} \partial_{\theta_1} \partial_{\theta_2} Q(\theta_2 | \theta_1) &= (-) \lim_{t \gg 0} \partial_{\theta_{t+1}}^2 \langle \log p(z | x, \theta_{t+1}) | x, \theta_t \rangle \\ &= (-) \langle \partial_{\theta_*}^2 \log p(z | x, \theta_*) | x, \theta_* \rangle = \mathcal{J}_{\text{cond}}(\theta_* | x)\end{aligned}$$

- Therefore, Rate of Convergence R of the parameter θ is proportional to ratios of Fisher matrices:

$$R = \lim_{t \gg 0} \frac{\theta_{t+1} - \theta_*}{\theta_t - \theta_*} = \frac{\mathcal{J}_{\text{com}}(\theta_* | x)}{\mathcal{J}_{\text{cond}}(\theta_* | x)}$$

Derivation

$$(\partial_{\theta} Q(\theta | \theta'))_{\theta_{(*)}} = 0 = (\partial_{\theta_{t+1}} Q(\theta_{t+1} | \theta_t)) \text{?}$$

- $\partial_{\theta_{t+1}} Q(\theta_t, \theta_{t+1}) = 0$ due to the definition of M Step:

$$\partial_{\theta_{t+1}} Q(\theta_t, \theta_{t+1}) = \lim_{\theta \rightarrow \theta_{t+1}} \partial_{\theta} \max_{\theta} Q(\theta | \theta_t) = 0$$

- θ_* maximizes $Q(\theta | \theta_*)$ for all θ :

$$(\partial_{\theta} Q(\theta | \theta'))_{\theta_{(*)}} = \lim_{\theta \rightarrow \theta_*} \partial_{\theta} \max_{\theta} Q(\theta | \theta_*) = 0$$

Derivation

$$"\mathcal{L}_N(\theta_t | x) \leq \mathcal{L}_N(\theta_{t+1} | x)"?$$

- Suppose decomposition: $\mathcal{L}_N(\theta_t | x) \stackrel{!}{=} F(x, \theta_t, q) + \Delta(x, \theta_t, q)$ where Δ is a non-negative term, with at least one zero $\Delta(q_*) = 0$.
- Suppose q_* sets $\Delta(x, \theta_t, q_*) = 0$:

$$\mathcal{L}_N(\theta_t | x) = F(x, \theta_t, q_*) \quad \text{"E Step"}$$

- Update $\theta_{t+1} \leftarrow \arg \max_{\theta_t} F(x, \theta_t, q_*)$ "M Step":

$$\mathcal{L}_N(\theta_t | x) \leq F(x, \theta_{t+1}, q_*) \leq \mathcal{L}_N(\theta_{t+1} | x) + (-)\Delta(x, \theta_{t+1}, q_*)$$

- Therefore, since Δ is non-negative:

$$\mathcal{L}_N(\theta_t | x) \leq \mathcal{L}_N(\theta_{t+1} | x)$$

Derivation

" $\Delta(x, \theta_t, q) = D_{\text{KL}}(q(z) \parallel p(z | x, \theta))$ "?

- Log-Likelihood fractures into various "entropies":

$$\begin{aligned}\mathcal{L}_N(\theta | x) &\equiv \sum_z q(z) \log p(x, \theta) = \sum_z q(z) \log (p(z, x, \theta) p(z | x, \theta)) \\&= \sum_z q(z) \log (p(z, x, \theta) p(z | x, \theta) q(z) q^{(-)}(z)) \\&= \sum_z q(z) \log q(z) + \sum_z q(z) \log (p(z, x, \theta)) \\&\quad + \sum_z q(z) \log (p(z | x, \theta) q^{(-)}(z)) \\&= (-)H(q) + Q(q(z) \parallel p(z, x, \theta)) + (-)D_{\text{KL}}(q(z) \parallel p(z | x, \theta))\end{aligned}$$

- Where $H(q) = (-) \sum_z q(z) \log q(z)$ is Shannon entropy. Assign:

$$\begin{aligned}F(x, \theta, q) &:= (-)H(q) + Q(q(z) \parallel p(z, x, \theta)) \\ \Delta(x, z, q) &:= D_{\text{KL}}(q(z) \parallel p(z | x, \theta))\end{aligned}$$

- Δ is non-negative since it's equivalent to Kullback–Leibler divergence:

$$\begin{aligned}\Delta(x, z, q) &= D_{\text{KL}}(q(z) \parallel p(z | x, \theta)) \\&= \sum_z q(z) \log q(z) p^{(-)}(z | x, \theta)\end{aligned}$$

- continued:

$$\sum_z q(z) \log q(z) p^{(-)}(z | x, \theta) = \left\langle \log q(z) p^{(-)}(z | x, \theta) \right\rangle_{z \sim q}$$

- Jensens' inequality for $\phi(\bullet) = (-)\log(\bullet)$ which is convex:

$$\left\langle \log q(z) p^{(-)}(z | x, \theta) \right\rangle_{z \sim q} = (-) \left\langle (-)\log q(z) p^{(-)}(z | x, \theta) \right\rangle = \left\langle \phi(q^{(-)}(z) p(z | x, \theta)) \right\rangle$$

- By that identity, KL divergences are non-negative:

$$\begin{aligned} \left\langle \phi(q^{(-)}(z) p(z | x, \theta)) \right\rangle &\geq \phi\left(\left\langle q^{(-)}(z) p(z | x, \theta) \right\rangle\right) \\ &\geq \log\left(\sum_z q(z) q^{(-)}(z) p(z | x, \theta)\right) \\ &\geq \log\left(\sum_z p(z | x, \theta)\right) \\ &\geq \log(1) \\ &\geq (0) \end{aligned}$$

- Therefore $\Delta(x, z, q) \geq (0)$
- To achieve the minimum: $q(z) \leftarrow p(z | x, \theta)$ "E Step":

$$\Delta(x, z, q) = \sum_z q(z) \log q(z) p^{(-)}(z | x, \theta) \stackrel{!}{=} 0 \quad \forall q(z) \neq 0$$

$$0 = \log\left(\sum_z q(z) p^{(-)}(z | x, \theta)\right)$$

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

summarize GMM:

- Log-likelihood of "hard clustering": $\mathcal{L}_N^{\text{GMM}} \equiv \sum_{\alpha \in [1, N]} \log p(x_\alpha | \theta)$
- Assume "latent" variables in "one-hot" representations: $z_\alpha \sim \mathbf{1}_\alpha(J_*)$ with $J_* \in [1, \mathcal{K}]$.
- z_α encodes cluster indices for data $x_\alpha, \forall \alpha \in [1, N]$. This means z_α distributes in accord with a multinomial distribution:

$$p(z | \pi) = \prod_{\alpha \in [1, N]} p(z_\alpha | \pi_I) = \prod_{\alpha \in [1, N]} \prod_{I \in [1, \mathcal{K}]} (\pi_I)^{z_{\alpha I}^I}$$

- This can be derived by substituting "one hot" $z_{\alpha I}$ into a general multinomial PD:

$$p(z_\alpha | \pi) = \frac{\mathcal{N}!}{\prod_{\alpha \in [1, N]} z_{\alpha I}} \prod_{\alpha \in [1, N]} (\pi_I)^{z_{\alpha I}} \quad \text{where} \quad \mathcal{N} := \sum_{\alpha \in [1, N]} z_{\alpha I}$$

- Given a label z_α , assume x_α follows a Gaussian distribution, and the joint PD:

$$p(x | z\theta) \equiv \prod_{\alpha} \prod_I \mathcal{N}(x_\alpha, (\mu\sigma)_{\alpha I})^{(z_\alpha)^I} := \prod_{\alpha} \prod_I g_{\alpha I}^{(z_\alpha)^I}$$

- Primitive attempt leads to "incomplete" likelihood:

$$\mathcal{L}_N(\theta | x) \equiv \log \prod_{\alpha \in [1, N]} p(x_\alpha, \theta) = \sum_{\alpha \in [1, N]} \log \sum_{z_\alpha} p(x_\alpha, z_\alpha, \theta)$$

GMM, continued:

- However, the task now involves "summing over all categories of z ," which is hard (Why?). Suppose z be chosen, $\sum_z \sim 1$ and it leads to "complete" likelihood:

$$\mathcal{L}_N^{\text{comp}}(\theta | x) = \sum_{\alpha \in [1, N]} \log p(x_\alpha, z_\alpha, \theta) = \sum_{\alpha \in [1, N]} \sum_{I \in [1, \mathcal{K}]} (z_\alpha)^I \log \pi_I \cdot g_{\alpha I}$$

- Define Q -function:

$$Q(\theta_{t+1} | \theta_t) \equiv \langle \mathcal{L}_N^{\text{comp}} | x \theta_t \rangle = \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} \left[\sum_{z_\alpha} p(z_\alpha | x \theta_t) (z_\alpha)^I \right] \log (\pi_I \cdot g_{\alpha I})$$

- The average over $(z_\alpha)^I$ is called Responsibility, and amounts to a probability weight for Q :

$$r_{\alpha I} \equiv \sum_{J \in [1, \mathcal{K}]} p(z_\alpha | x \theta^t) (z_\alpha)^J = p(z_\alpha^I | x \theta^t) = \frac{\pi_I g_{I\alpha}}{\sum_{J \in [1, \mathcal{K}]} \pi_J g_{J\alpha}}$$

- Thus, Q -function works out to be:

$$Q(\theta_{t+1} | \theta_t) \equiv \langle \mathcal{L}_N^{\text{comp}} | x \theta_t \rangle = \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha I} \cdot \log (\pi_I \cdot g_{\alpha I})$$

- EM algorithm optimizes Q -function on θ :

- E Step consists of evaluating $r_{I\alpha}$ and then $Q(\theta^{t+1} | \theta^t)$.
- M Step consists of updating the parameters θ by max'ing out Q :

$$\theta_I^{t+1} \leftarrow \arg \max_{\theta} Q(\theta | \theta^t)$$

GMM, continued:

- The update rules work out to be:

$$\pi_I^{t+1} \leftarrow \frac{N_I}{N} \quad N_I \equiv \sum_{i \in [1, N]} r_{\alpha I}^t$$
$$\mu_I^{t+1} \leftarrow \frac{1}{N_I} \left(\sum_{\alpha \in [1, N]} r_{I\alpha}^t x_\alpha \right) \quad \sigma_I^{t+1} \leftarrow \frac{\sum_{\alpha} r_{\alpha I} \|x_\alpha - \mu_I\|^2}{2 \sum_{\alpha} r_{\alpha I}}$$

Derivation

$$\pi_I^{t+1} \leftarrow \frac{N_I}{N} \quad N_I \equiv \sum_{\alpha \in [1, N]} r_{\alpha I}^t ?$$

- Remind that the task is to update π as the maximal solution to Q :

$$\pi_I^{t+1} \leftarrow \arg \max_{\pi_I} \left[Q(\theta | \theta^t) + \lambda \cdot \left(\sum_{I \in [1, \mathcal{K}]} \pi_I - 1 \right) \right]$$

- In practice, differentiate Q with the Lagrange multiplier λ , and set to 0:

$$\begin{aligned} 0 &\stackrel{!}{=} \partial_{\pi_I} \left(Q(\theta | \theta^t) + \lambda \cdot \left(\sum_{J \in [1, \mathcal{K}]} \pi_J - 1 \right) \right) = \partial_{\pi_I} Q(\theta | \theta^t) + \lambda \sum_{J \in [1, \mathcal{K}]} \partial_{\pi_I} \pi_J \\ &= \partial_{\pi_I} Q(\theta | \theta^t) + \lambda \sum_{J \in [1, \mathcal{K}]} \delta_{IJ} = \partial_{\pi_I} Q(\theta | \theta^t) + \lambda \end{aligned}$$

- The derivative of Q :

$$\begin{aligned} \partial_{\pi_I} Q(\theta | \theta^t) &= \partial_{\pi_I} \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha I} \log(\pi_I \cdot g_{\alpha I}) = \sum_{I \in [1, \mathcal{K}]} r_{\alpha I} \partial_{\pi_I} \log(\pi_I \cdot g_{\alpha I}) \sum_{I \in [1, \mathcal{K}]} r_{\alpha I} (\pi_I) \\ &= \sum_{I \in [1, \mathcal{K}]} \pi_I^{(-)} r_{\alpha I} \delta_{IJ} = \sum_{\alpha} \pi_J^{(-)} r_{\alpha J} = \pi_J^{(-)} N_J \end{aligned}$$

- Solving the differential equation:

$$0 \stackrel{!}{=} \partial_{\pi_I} Q(\theta | \theta^t) + \lambda = \pi_J^{(-)} N_J + \lambda \implies \pi_J = (-) \frac{N_J}{\lambda}$$

- To obtain λ , trace over that identity in the last item \sum_J :

$$\begin{aligned}
 (1) &= \sum_{J \in [1, \mathcal{K}]} \pi_J = \sum_{J \in [1, \mathcal{K}]} (-) \frac{N_J}{\lambda} = (-) \frac{1}{\lambda} \sum_{J \in [1, \mathcal{K}]} \left(\sum_{\alpha \in [1, N]} r_{\alpha J} \right) \\
 &= (-) \frac{1}{\lambda} \sum_{\alpha \in [1, N]} \left(\sum_{J \in [1, \mathcal{K}]} r_{\alpha J} \right) = (-) \frac{1}{\lambda} \sum_{\alpha \in [1, N]} (1) = (-) \frac{1}{\lambda} (N)
 \end{aligned}$$

- WLOG, notice $r_{\alpha J}$ is the posterior PD, which is normalized to 1:

$$\sum_{J \in [1, \mathcal{K}]} r_{\alpha J} = \sum_{J \in [1, \mathcal{K}]} p(z_{\alpha}^J | x, \theta) \stackrel{!}{=} 1$$

- Therefore, the update rule for π_J :

$$\lambda = (-)N \implies \pi_J = (-) \frac{N_J}{\lambda} = \frac{N_J}{N}$$

Derivation

$$\mu_I^{t+1} \leftarrow \frac{1}{N_I} (\sum_{\alpha \in [1, N]} r_{I\alpha}^t x_\alpha) \text{?}$$

- Similarly, the update rule for μ_I can be obtained taking a derivative of Q over it:

$$\begin{aligned} \partial_{\mu_I} Q(\theta | \theta^t) &= \partial_{\mu_I} \sum_{J \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha J} \log(\pi_J \cdot g_{\alpha J}) \\ &= \sum_{J \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha J} \partial_{\mu_I} \log(\pi_J \cdot g_{\alpha J}) \end{aligned}$$

- evaluate the log content first, and then evaluate the derivative:

$$\begin{aligned} \partial_{\mu_I} \log(\pi_J \cdot g_{\alpha J}) &= \partial_{\mu_I} \log \left[\mathcal{Z}_J e^{(-)2^{(-)} \text{Tr}((x_\alpha - \mu_J)^\top \sigma_J^{(-)}(x_\alpha - \mu_J))} \right] \\ &= \partial_{\mu_I} \left[\log \mathcal{Z}_J + (-)2^{(-)} \text{Tr}((x_\alpha - \mu_J)^\top \sigma_J^{(-)}(x_\alpha - \mu_J)) \right] \\ &= \partial_{\mu_I} \left[(2^{(-)} \dim x_\alpha) \log(2\pi \|\sigma_J\|)^{(-)} \right. \\ &\quad \left. + (-)2^{(-)} \text{Tr}((x_\alpha - \mu_J)^\top \sigma_J^{(-)}(x_\alpha - \mu_J)) \right] \\ &= \partial_{\mu_I} \left[(0) + (-)2^{(-)} \text{Tr}((x_\alpha - \mu_J)^\top \sigma_J^{(-)}(x_\alpha - \mu_J)) \right] \\ &= (-)2^{(-)} \left((x_{\alpha m} - \mu_{Jm}) [\sigma_J^{(-)}]^{mn} (x_{\alpha n} - \mu_{Jn}) \right) \\ &= (x_{\alpha m} - \mu_{Jm}) [\sigma_J^{(-)}]^{mn} (\delta_{IJ} \mathbf{1}_n) \\ &= \delta_{IJ} \cdot (x_\alpha - \mu_J)^\top \sigma_J^{(-)} \mathbf{1} \end{aligned}$$

- Thus, substitute and evaluate the derivative of Q :

$$\begin{aligned}
 \partial_{\mu_I} Q(\theta | \theta^t) &= \sum_{J \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha J} \partial_{\mu_I} \log(\pi_J \cdot g_{\alpha J}) \\
 &= \sum_{J \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha J} \delta_{IJ} \cdot (x_\alpha - \mu_J)^\top \sigma_J^{(-)} \mathbf{1} \\
 &= \sum_{\alpha \in [1, N]} r_{\alpha I} (x_\alpha - \mu_I)^\top \sigma_I^{(-)} \mathbf{1}
 \end{aligned}$$

- setting to 0 and invert the differential equation:

$$\begin{aligned}
 0 \stackrel{!}{=} \partial_{\mu_I} Q(\theta | \theta^t) &= \sum_{\alpha \in [1, N]} r_{\alpha J} (x_\alpha - \mu_J)^\top \sigma_J^{(-)} \mathbf{1} \\
 &= \sum_{\alpha \in [1, N]} r_{\alpha J} x_\alpha^\top \sigma_J^{(-)} \mathbf{1} + (-) \sum_{\alpha \in [1, N]} r_{\alpha J} \mu_J^\top \sigma_J^{(-)} \mathbf{1}
 \end{aligned}$$

- Therefore, the update rule for μ_J :

$$\begin{aligned}
 \sum_{\alpha \in [1, N]} r_{\alpha J} x_\alpha^\top \sigma_J^{(-)} \mathbf{1} &= \sum_{\alpha \in [1, N]} r_{\alpha J} \mu_J^\top \sigma_J^{(-)} \mathbf{1} \\
 \mu_J^{t+1} &\leftarrow \frac{\sum_{\alpha \in [1, N]} r_{\alpha J} x_\alpha}{\sum_{\alpha \in [1, N]} r_{\alpha J}}
 \end{aligned}$$

Derivation

$$\|\sigma_I^{t+1} \leftarrow \frac{\sum_{\alpha} r_{\alpha I} \|x_{\alpha} - \mu_I\|^2}{2 \sum_{\alpha} r_{\alpha I}}\|?$$

- For ease of notation, denote $\partial_{\sigma_I^{(-)}} \rightarrow \partial_I$
- As a rhyme to the two derivations above, differentiate Q over $\sigma_I^{(-)}$:

$$\begin{aligned} \partial_I Q(\theta | \theta^t) &= \sum_{\alpha, I} r_{\alpha I} \partial_I \left[\log \mathcal{Z}_J^{(-)} + (-) \frac{1}{2} \text{Tr}((x_{\alpha} - \mu_J)^{\top} \sigma_J^{(-)} (x_{\alpha} - \mu_J)) \right] \\ &= \sum_{\alpha, I} r_{\alpha I} \left[\partial_I \log \mathcal{Z}_J^{(-)} + (-) \frac{1}{2} \partial_I \text{Tr}((x_{\alpha} - \mu_J)^{\top} \sigma_J^{(-)} (x_{\alpha} - \mu_J)) \right] \end{aligned}$$

- Applying a well known matrix identity, i.e $\frac{\delta \det g}{\det g} = \delta g g^{(-)}$

$$\begin{aligned} \partial_I \log \mathcal{Z}_J^{(-)} &= \partial_I \log (2\pi \|\sigma_J\|)^{(-)} = \partial_I \log ((2\pi)^{(-)} \det \sigma_J^{(-)}) = \partial_I \log ((1) \det \sigma_J^{(-)}) \\ &= \partial_I \log (\det \sigma_J^{(-)})^{(-)} = \sigma_J \delta_{IJ} \mathbf{1}_{\zeta \times \zeta} \quad \text{where} \quad \zeta := \dim x_{\alpha} \end{aligned}$$

- Cyclic permute the trace content, and apply $\partial_A \text{Tr}(AB) = B$:

$$\begin{aligned} (-) \frac{1}{2} \partial_I \text{Tr}((x_{\alpha} - \mu_J)^{\top} \sigma_J^{(-)} (x_{\alpha} - \mu_J)) &= (-) \frac{1}{2} \partial_I \text{Tr}(\sigma_J^{(-)} (x_{\alpha} - \mu_J) (x_{\alpha} - \mu_J)^{\top}) \\ &= (-) \frac{1}{2} \text{Tr}(\partial_I \sigma_J^{(-)} (x_{\alpha} - \mu_J) (x_{\alpha} - \mu_J)^{\top}) \\ &= (-) \frac{1}{2} \text{Tr}(\delta_{IJ} \mathbf{1}_{\zeta \times \zeta} (x_{\alpha} - \mu_J) (x_{\alpha} - \mu_J)^{\top}) \\ &= (-) \frac{1}{2} \delta_{IJ} \|x_{\alpha} - \mu_J\|^2 \end{aligned}$$

- Finish evaluating the differentiation on Q , set to 0:

$$\begin{aligned}\partial_I Q(\theta | \theta^t) &= \sum_{\alpha, I} r_{\alpha I} \partial_I \log \mathcal{Z}_J + \sum_{\alpha, I} r_{\alpha I} (-) \frac{1}{2} \partial_I \text{Tr}((x_\alpha - \mu_J)^\top \sigma_J^{(-)} (x_\alpha - \mu_J)) \\ &= \sum_{\alpha} r_{\alpha I} \sigma_I \mathbf{1}_{\zeta \times \zeta} + \sum_{\alpha} r_{\alpha I} (-) \frac{1}{2} \|x_\alpha - \mu_I\|^2 \stackrel{!}{=} 0\end{aligned}$$

- Invert the differential equation results in:

$$\frac{1}{2} \sum_{\alpha} r_{\alpha I} \|x_\alpha - \mu_I\|^2 = \sum_{\alpha} r_{\alpha I} \sigma_I$$

- Thus, the update rule on σ_I :

$$\sigma_I^{t+1} \leftarrow \frac{\sum_{\alpha} r_{\alpha I} \|x_\alpha - \mu_I\|^2}{2 \sum_{\alpha} r_{\alpha I}}$$

K-means

- Duality K-means \sim GMM: Introduce generalized log-likelihood
 $\mathcal{L}_\epsilon(x^{\otimes N}, I) \equiv \epsilon^{(-)} \cdot \sum_{\alpha \in (1, \dots, N)} \log \sum_{I \in (1, \dots, \mathcal{K})} p(x_\alpha, I)^\epsilon$
 - $\lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) \approx (-)2^{(-)} \sum_{J \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N_J]} \|x_\alpha - \mu_J\|^2$
 - Both GMM and K-Means use EM algorithms. (i.e reducing to K-Means by substituting $\sigma_J \sim 1$).
- For *soft clustering* (def.?), introduce tight breaker r as an indicator function:

$$r_{\alpha I} = \mathbf{1}(\|x_\alpha - \mu_I\|^2 \leq \|x_\alpha - \mu_J\|^2) \quad \forall J \in [1, \mathcal{K}]$$

$$\mathcal{L}_N(x, I) := \sum_{\alpha \in [1, N]} \sum_{J \in [1, \mathcal{K}]} r_{\alpha J} \|x_\alpha - \mu_J\|^2$$

- algorithm:
 - assign clusters $\{c_I\} \leftarrow \arg \min \|x_\alpha - \mu_I\|^2 \quad \forall \alpha \in [1, N]$
 - $\mu_I \leftarrow \arg \max \mathcal{L}_N(x^{\otimes N}, I) \leftarrow (\sum_{\beta \in [1, N]} r_{\beta I})^{(-)} \cdot \sum_{\alpha \in [1, N]} r_{\alpha I} x_I$
- Aside: Suppose sample average \bar{x} : $\arg \max_{\mu_I} \mathcal{L}_N(x, I) \approx \arg \min_{\mu_I} B(x)$
 i.e $\sum_{\alpha, I} r_{\alpha I} W_I(x) + N \cdot B(x) \propto N^2 T(x) + \mathcal{O}(x, \mu, \bar{x}) \quad T \perp \{\mu_I\}$

$$T = N^{(-)} \sum_{\alpha} \|x_\alpha - \bar{x}\| \quad \text{total deviation}$$

$$W_I = (\sum_{\alpha} r_{\alpha I} \|x_\alpha - \mu_I\|^2) (\sum_{\beta} r_{\beta I})^{(-)} \quad \text{intra-cluster deviation}$$

$$B = \sum_I N^{(-)} \sum_{\alpha} r_{\alpha I} \|\mu_I - \bar{x}\|^2 \quad \text{inter-cluster deviation}$$

Derivation

$$\lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) = \sum_{J, \alpha_J} \log p(J, x_{\alpha_J})$$

- Note: $\mathcal{L}_\epsilon(x^{\otimes N}, I)$ at $\epsilon \rightarrow 1$ corresponds to GMM:

$$\lim_{\epsilon \rightarrow 1} \mathcal{L}_\epsilon(x^{\otimes N}, I) \equiv \lim_{\epsilon \rightarrow 1} \epsilon^{(-)} \cdot \sum_{\alpha \in [1, N]} \log \sum_{I \in [1, \mathcal{K}]} p(x_\alpha, I)^\epsilon \approx \sum_{\alpha \in [1, N]} \log \sum_{I \in [1, \mathcal{K}]} p(x_\alpha, I)$$

- The promotion $\mathcal{L}(I, x_\alpha) \mapsto \mathcal{L}_\epsilon(I, x_\alpha)$ requires $p(I | x_\alpha) \mapsto p_\epsilon(I | x_\alpha)$:

$$p(I | x_\alpha) = (p(x_\alpha, I))(p(x_\alpha))^{(-)} = (p(x_\alpha, I)) \left(\sum_{J \in [1, \mathcal{K}]} p(x_\alpha, J) \right)^{(-)}$$

$$p_\epsilon(I | x_\alpha) \equiv (p^\epsilon(x_\alpha, I)) \left(\sum_{J \in [1, \mathcal{K}]} p^\epsilon(x_\alpha, J) \right)^{(-)}$$

- Assume there is a cluster I so that $p(x_\alpha, I) > p(x_\alpha, J)$ for all $J \neq I, J \in [1, \mathcal{K}]$ "E Step" (Why?)

$$\begin{aligned} \lim_{\epsilon \gg 0} p_\epsilon^{(-)}(I | x_\alpha) &= \lim_{\epsilon \gg 0} (p^\epsilon(x_\alpha, I))^{(-)} \left(\sum_{J \in [1, \mathcal{K}]} p^\epsilon(x_\alpha, J) \right) = \lim_{\epsilon \gg 0} \sum_{J \in [1, \mathcal{K}]} \left(\frac{p(x_\alpha, J)}{p(x_\alpha, I)} \right)^\epsilon \\ &= \lim_{\epsilon \gg 0} \left(1 + \sum_{J \neq I} \left(\frac{p(x_\alpha, J)}{p(x_\alpha, I)} \right)^\epsilon \right) = 1 + \sum_{J \neq I} \lim_{\epsilon \gg 0} \left(\frac{p(x_\alpha, J)}{p(x_\alpha, I)} \right)^\epsilon \\ &= 1 + \lim_{\epsilon \gg 0} \sum_{J \neq I} (0) = 1 \end{aligned}$$

- The following items assume some simple limits:

$$\lim_{\epsilon \gg 0} p_\epsilon(I | x_\alpha) = 1 \quad \lim_{\epsilon \gg 0} p_\epsilon(I | x_\alpha) \log p_\epsilon(I | x_\alpha) = 1 \log 1 = 0$$

- In the large ϵ limit, introducing a (1) in the generalized likelihood:

$$\begin{aligned} \lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) &= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{\alpha \in [1, N]} \log \sum_{I \in [1, \mathcal{K}]} p(x_\alpha, I)^\epsilon \\ &= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{\alpha \in [1, N]} \left(\sum_{J \in [1, \mathcal{K}]} p_\epsilon(J | x_\alpha) \right) \log \sum_{I \in [1, \mathcal{K}]} p^\epsilon(x_\alpha, I) \\ &= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{\alpha \in [1, N]} \left(\sum_{J \in [1, \mathcal{K}]} p_\epsilon(J | x_\alpha) \right) \log \frac{p_\epsilon(J, x_\alpha)}{p^\epsilon(J | x_\alpha)} \\ &= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{\alpha \in [1, N]} \sum_{J \in [1, \mathcal{K}]} p_\epsilon(J | x_\alpha) \log \frac{p_\epsilon(J, x_\alpha)}{p^\epsilon(J | x_\alpha)} \end{aligned}$$

- If for each $J \in \mathcal{K}$, $N(J)$ is not uniform, i.e $N(J) = N_J$, the double sum is replaced by:

$$\sum_{\alpha \in [1, N]} \sum_{J \in [1, \mathcal{K}]} \rightarrow \sum_{J \in [1, \mathcal{K}]} \sum_{\alpha_J \in [1, N_J]} = \sum_{J, \alpha_J}$$

- Therefore, with ϵ being very large, the likelihood limits to:

$$\lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) = \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{J, \alpha_J} p_\epsilon(J | x_{\alpha_J}) \log \frac{p_\epsilon(J, x_{\alpha_J})}{p^\epsilon(J | x_{\alpha_J})}$$

- continued:

$$\begin{aligned}
\lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) &= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{J, \alpha_J} p_\epsilon(J | x_{\alpha_J}) \log p_\epsilon(J, x_{\alpha_J}) \\
&\quad + (-) \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{J, \alpha_J} p_\epsilon(J | x_{\alpha_J}) \log p^\epsilon(J | x_{\alpha_J}) \\
&= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{J, \alpha_J} (1) \log p^\epsilon(J, x_{\alpha_J}) + (-) \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{J, \alpha_J} (0) \\
&= \lim_{\epsilon \gg 0} \epsilon^{(-)} \sum_{J, \alpha_J} \log p^\epsilon(J, x_{\alpha_J}) \\
&= \lim_{\epsilon \gg 0} \sum_{J, \alpha_J} \log p^{\epsilon \cdot \epsilon^{(-)}}(J, x_{\alpha_J}) \\
&= \sum_{J, \alpha_J} \log p(J, x_{\alpha_J})
\end{aligned}$$

Derivation

$$\|\mathcal{L}_{\epsilon \gg 0}^{\text{GMM}}(I, \theta | x^{\otimes N}) = (-)2^{(-)} \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha_I \in [1, N_I]} \|x_{\alpha_I} + (-)\mu_I\|^2\|?$$

- Factorize the joint PD by Bayes' rule, and substitute in Normal Mixtures (GMM) WLOG:

$$p(J, x_{\alpha_J}) = p(J | x_{\alpha_J}, \theta) p(x_{\alpha_J} | \theta) = \mathcal{Z}_J^{(-)} e^{(-)\frac{1}{2} \text{Tr}(x_{\alpha_J} - \mu_J)^\top \sigma_J^{(-)} (x_{\alpha_J} - \mu_J)} \pi_J$$

- Substitute the aforementioned item into the likelihood:

$$\begin{aligned} \lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) &= \sum_{J, \alpha_J} \log \mathcal{Z}_J^{(-)} e^{(-)\frac{1}{2} \text{Tr}(x_{\alpha_J} - \mu_J)^\top \sigma_J^{(-)} (x_{\alpha_J} - \mu_J)} \pi_J \\ &= \sum_{J, \alpha_J} \log \mathcal{Z}_J^{(-)} + \sum_{J, \alpha_J} (-)\frac{1}{2} \text{Tr}(x_{\alpha_J} - \mu_J)^\top \sigma_J^{(-)} (x_{\alpha_J} - \mu_J) \end{aligned}$$

- The first term above:

$$\begin{aligned} \mathcal{Z}_J(x_{\alpha_J})^{(-)} &= ((2\pi)^\zeta \det\{\sigma_J\})^{(-)2^{(-)}} \quad \zeta := \dim x_{\alpha_J} \\ \sum_{J, \alpha_J} \log \mathcal{Z}_J^{(-)} &= \sum_{J, \alpha_J} \log ((2\pi)^\zeta \det\{\sigma_J\})^{(-)2^{(-)}} = (-)\frac{1}{2} \sum_{J, \alpha_J} \log ((2\pi)^\zeta \det\{\sigma_J\}) \end{aligned}$$

- Assume that $\sigma_J \stackrel{!}{=} \mathbf{1}_{\zeta \times \zeta}$
- Assume also that $\mathcal{Z}_J^{(-)}(x_{\alpha_J}) \sim (0)$ (WLOG?)

- continued:

$$\begin{aligned}\lim_{\epsilon \gg 0} \mathcal{L}_\epsilon(x^{\otimes N}, I) &\propto (0) + \sum_{J, \alpha_J} (-) \frac{1}{2} \text{Tr}(x_{\alpha_J} - \mu_J)^\top (1)(x_{\alpha_J} - \mu_J) \\ &\propto (-) \frac{1}{2} \sum_{J, \alpha_J} \|x_{\alpha_J} - \mu_J\|^2\end{aligned}$$

Derivation

$$\mu_I^{t+1} \leftarrow \sum_{\alpha_I} x_{\alpha_I} \text{?}$$

- Only M Step requires derivation.
- For *hard clustering*, $r_{\alpha_I I} \equiv 1$. π_I are stationary:

$$\sum_{\alpha \in [1, N_I]} r_{\alpha_I I} = \sum_{\alpha \in [1, N_I]} (1) = N_I$$

- WLOG: $\sigma_I(x_{\alpha_I}) \sim 1$. The remaining update:

$$\mu_I^{t+1} \leftarrow \arg \max_{\mu_I} \mathcal{L}_{\epsilon \sim 1}(x^{\otimes N}, I) = N_I^{(-)} \cdot \sum_{i \in [1, N_I]} x_{\alpha_I}$$

Derivation

" $\mu_I \leftarrow N_I^{(-)} \cdot \sum_{\alpha \in [1, \mathcal{K}]} r_{\alpha I} x_I$ "? where $N_I \equiv \sum_{\beta \in [1, \mathcal{K}]} r_{\beta I}$ "?

- Assume the task: $\mu_I \leftarrow \arg \max_{\mu_I} \mathcal{L}_N(x^{\otimes N}, I)$:

$$\begin{aligned} 0 &\stackrel{!}{=} \partial_{\mu_I} \mathcal{L}_N(x^{\otimes N}, I) = \partial_{\mu_I} \sum_{\alpha \in [1, N]} r_{\alpha I} \|x_{\alpha} - \mu_I\|^2 \\ &= \sum_{\alpha \in [1, N]} r_{\alpha I} (x_{\alpha} - \mu_I)(-2) \end{aligned}$$

- Therefore:

$$\mu_I := \frac{\sum_{\alpha \in [1, N]} r_{\alpha I} x_{\alpha}}{\sum_{\beta \in [1, N]} r_{\beta I}}$$

Derivation

$$"\sum_{\alpha, I} r_{\alpha I} W_I(x) + N \cdot B(x) \propto N^2 T(x) + \mathcal{O}(x, \mu, \bar{x}) \quad T \perp \{\mu_I\}"?$$

- The first piece on the left is equivalent to the Score function.

$$\begin{aligned} \sum_{\alpha, I} r_{\alpha I} W_I(x) &= \sum_{\alpha \in [1, N]} \sum_{I \in [1, \mathcal{K}]} r_{\alpha I} \left(\sum_{\gamma \in [1, N]} r_{\gamma I} \|x_{\gamma} - \mu_I\|^2 \right) \left(\sum_{\beta} r_{\beta I} \right)^{(-)} \\ &= \sum_{I \in [1, \mathcal{K}]} \sum_{\gamma \in [1, N]} r_{\gamma I} \left(\sum_{\alpha \in [1, N]} r_{\alpha I} \|x_{\alpha} - \mu_I\|^2 \right) \left(\sum_{\beta \in [1, N]} r_{\beta I} \right)^{(-)} \\ &= \sum_{I \in [1, \mathcal{K}]} \sum_{\gamma \in [1, N]} r_{\gamma I} \|x_{\gamma} - \mu_I\|^2 = \mathcal{L}_N(x, \mu) \end{aligned}$$

- A term proportional to T can be factored out from the Left Hand Side:

$$\begin{aligned} \sum_{\alpha, I} r_{\alpha I} W_I(x) + N \cdot B(x) &= \sum_{I \in [1, \mathcal{K}]} \sum_{\gamma \in [1, N]} r_{\gamma I} \|x_{\gamma} - \mu_I\|^2 + \frac{N}{N} \cdot \sum_{I, \alpha} r_{\alpha I} \|\mu_I - \bar{x}\|^2 \\ &= \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha I} (\|x_{\alpha} - \mu_I\|^2 + \|\mu_I - \bar{x}\|^2) \\ &= \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha I} (x_{\alpha}^2 + (2)\mu_I^2 + \bar{x}^2 + (-2)x_{\alpha}\mu_I + (-2)\mu_I\bar{x}) \\ &= \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha I} (\|x_{\alpha} + (-)\bar{x}\|^2 + (-2)x_{\alpha}\mu_I + (2)\mu_I^2 + x_{\alpha}) \\ &= \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{\alpha I} \|x_{\alpha} + (-)\bar{x}\|^2 + \mathcal{O}(x_{\alpha}, \bar{x}, \mu_I) \end{aligned}$$

- WLOG: " $\sum_{I \in [1, \mathcal{K}]} r_{I\alpha} \stackrel{?}{=} N \quad \forall r \propto N \times N$ "

$$\begin{aligned} \sum_{\alpha, I} r_{\alpha I} W_I(x) + N \cdot B(x) &= \sum_{\alpha \in [1, N]} (N) \|x_\alpha + (-)\bar{x}\|^2 + \mathcal{O}(x_\alpha, \bar{x}, \mu_I) \\ &= (N)^2 \cdot T(x) + \mathcal{O}(x_\alpha, \bar{x}, \mu_I) \end{aligned}$$

- Where $\mathcal{O}(x_\alpha, \bar{x}, \mu_I) \equiv \sum_{I \in [1, \mathcal{K}]} \sum_{\alpha \in [1, N]} r_{I, \alpha} (x_\alpha + (-2)x_\alpha \mu_I + (2)\mu_I^2)$
- "Bias-Variance Decomp?" in disguise?

Derivation

$$p(z_\alpha^I | x\theta^t) = \frac{\pi_I g_{I\alpha}}{\sum_{J \in [1, \mathcal{K}]} \pi_J g_{J\alpha}} \text{?}$$

- Factorize x into x_β , and note $p(z_\alpha^I | x_\beta \theta^t) \stackrel{!}{=} 0$ if $\alpha \neq \beta$:

$$p(z_\alpha^I | x\theta^t) = \prod_{\beta \in [1, N]} p(z_\alpha^I | x_\beta \theta^t) = \prod_{\beta \in [1, N]} p(z_\alpha^I | x_\alpha \theta^t)^{\delta_{\alpha\beta}} = p(z_\alpha^I | x_\alpha \theta^t)$$

- Apply Bayes rule:

$$p(z_\alpha^I | x_\alpha \theta^t) p(x_\alpha | \theta^t) = p(z_\alpha^I, x_\alpha | \theta^t) = p(x_\alpha | z_\alpha^I, \theta^t) p(z_\alpha^I | \theta^t) = \pi_I g_{\alpha I}$$

- solving for the posterior:

$$p(z_\alpha^I | x_\alpha \theta^t) = \frac{\pi_I g_{\alpha I}}{p(x_\alpha | \theta^t)}$$

- For $p(x_\alpha | \theta^t)$, introduce $\{z_\alpha^J\}$:

$$p(x_\alpha | \theta^t) = \sum_{J \in [1, \mathcal{K}]} p(x_\alpha z_\alpha^J | \theta^t) = \sum_{J \in [1, \mathcal{K}]} \pi_J g_{\alpha J}$$

- Therefore:

$$\therefore p(z_\alpha^I | x\theta^t) = \frac{\pi_I g_{I\alpha}}{\sum_{J \in [1, \mathcal{K}]} \pi_J g_{J\alpha}}$$

Silhouette, AIC, & BIC

- Denote *subcluster* to be \mathcal{C}_I . For all points $(x_{\alpha_I}, y_{\alpha_I}) \in \cup_{I \in [1, \mathcal{K}]} \mathcal{C}_I$, **Silhouette Coefficient** S is the **net discrepancy** (mod normalization) between the intra-cluster mean distances and the *minimal* inter-cluster mean distances:

$$S \equiv \sum_{I=1}^{\mathcal{K}} \sum_{\alpha_I=1}^{\mathcal{N}_I} \hat{s}(y_{\alpha_I}) \quad \text{where} \quad \hat{s}(y_{\alpha_I}) \equiv \frac{\hat{b}_{\alpha_I} - \hat{a}_{\alpha_I}}{\max(\hat{a}_{\alpha_I}, \hat{b}_{\alpha_I})}$$
$$\hat{b}_{\alpha_I} \equiv \min_{J \in \mathcal{C}, J \neq I} \sum_{\beta_J=1}^{\mathcal{N}_J} \frac{d(\alpha_I, \beta_J)}{\mathcal{N}_I - 1} \quad \hat{a}_{\alpha_I} \equiv \sum_{\substack{\beta_I=1 \\ \beta_I \neq \alpha_I}}^{\mathcal{N}_I} \frac{d(\alpha_I, \beta_I)}{\mathcal{N}_I}$$

- AIC:** For $X \equiv \{X_{\alpha}\}_{\alpha \in [1, \mathcal{N}]}$ random IID, AIC is the **unbiased** estimator of the **true risk** of the log-likelihood $\mathcal{L}_{\mathcal{N}}$ (mod a factor $\times 2$):

$$\text{AIC} = 2\mathcal{K} - 2 \max_{\theta} \mathcal{L}_{\mathcal{N}}$$

- BIC:** For $X \equiv \{X_{\alpha}\}_{\alpha \in [1, \mathcal{N}]}$ random IID, BIC is the first order approximation ($\sim \mathcal{O}(\mathcal{N})$) of the **log evidence** $\log(X | m)$ for all models $m_I, \forall I \in [1, \mathcal{K}]$:

$$\text{BIC} = \mathcal{K} \log \mathcal{N} - \max_{\theta} \mathcal{L}_{\mathcal{N}}$$

Derivation

$$\text{"AIC} \stackrel{?}{=} 2\mathcal{H} - 2\mathcal{L}_{\mathcal{N}}"$$

- Review MLE
- Consider log-likelihood \mathcal{L} and mean log-likelihood $\overline{\mathcal{L}}$

$$\mathcal{L}(\theta | x) \equiv \log p(x | \theta) = \log \prod_{\alpha=1}^{\mathcal{N}} p(x_{\alpha} | \theta) = \sum_{\alpha=1}^{\mathcal{N}} \log p(x_{\alpha} | \theta)$$

$$\overline{\mathcal{L}}(\theta | x) \equiv \langle \mathcal{L}(\theta | x) \rangle_{x \sim p} = \int dp(x) \mathcal{L}(\theta | x)$$

- goals of MLE are to optimize \mathcal{L} and $\overline{\mathcal{L}}$

$$\hat{\theta} \equiv \arg \max_{\theta} \mathcal{L} \quad \theta^* \equiv \arg \max_{\theta} \overline{\mathcal{L}} \quad \text{i.e.} \quad \max_{\theta} \mathcal{L}(\theta | x) = \mathcal{L}(\hat{\theta})$$

- Introduce True Risk $\mathfrak{R}(\mathcal{L})$ and Empirical Risk $\hat{\mathfrak{R}}(\mathcal{L})$

$$\mathfrak{R}(\mathcal{L}(\hat{\theta} | x)) \equiv [-]_{\mathcal{N}} \overline{\mathcal{L}}(\hat{\theta} | x) \quad \hat{\mathfrak{R}}(\mathcal{L}(\hat{\theta} | x)) \equiv \sum_{\alpha=1}^{\mathcal{N}} [-] \mathcal{L}(\hat{\theta} | x_{\alpha})$$

- $\hat{\mathfrak{R}}$ serves as the estimator for \mathfrak{R}
- AIC is $\hat{\mathfrak{R}}$ mod bias. To find the bias of \mathfrak{R}

$$\text{bias}(\mathfrak{R}) \equiv \mathbb{E}(\hat{\mathfrak{R}}) - \mathfrak{R}$$

- Approximate $\mathcal{L}(\hat{\theta})$ near the fixed point θ^* via Taylor Expansion

$$\mathcal{L}(\hat{\theta} | x) \equiv \sum_{\alpha=1}^{\mathcal{N}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \quad \mathcal{L}(\hat{\theta} | x_{\alpha}) = \log p(\hat{\theta} | x_{\alpha})$$

Derivation

- continued. Taylor expand around $\hat{\theta} \sim \theta^*$

$$\begin{aligned}\mathcal{L}(\hat{\theta} | x) &\approx \sum_{\alpha=1}^{\mathcal{N}} \mathcal{L}(\theta^* | x_{\alpha}) + (\hat{\theta} - \theta^*) \sum_{\alpha=1}^{\mathcal{N}} \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}}^{\top} \\ &\quad + \frac{1}{2} (\hat{\theta} - \theta^*)^{\top} \left[\sum_{\alpha=1}^{\mathcal{N}} \nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \mathcal{L}(\theta^* | x) \right]_{\theta^*=\hat{\theta}} (\hat{\theta} - \theta^*)\end{aligned}$$

- Used below (since $\hat{\theta} \equiv \arg \max_{\theta} \mathcal{L}(\theta | x)$)

$$0 = \nabla_{\hat{\theta}} \mathcal{L}(\theta | x) = \nabla_{\hat{\theta}} \log \prod_{\alpha=1}^{\mathcal{N}} p(\theta | x_{\alpha}) = \sum_{\alpha=1}^{\mathcal{N}} \nabla_{\hat{\theta}} \mathcal{L}(\theta | x_{\alpha})$$

- Consider $\mathcal{O}(\theta)$ term

$$\begin{aligned}&\sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}}^{\top} \\ &= \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}}^{\top} - 0 \\ &= \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}}^{\top} - \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \mathcal{L}(\theta | x_{\alpha}) \right]^{\top}\end{aligned}$$

Derivation

- continued.

$$\begin{aligned}
 & \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}}^{\top} \\
 &= \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left\{ \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} - \nabla_{\hat{\theta}} \mathcal{L}(\theta^* | x_{\alpha}) \right\}^{\top} \\
 &\approx \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*) \left\{ \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} - (\hat{\theta} - \theta^*)^{\top} \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} \right. \\
 &\quad \left. - \left[\nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} \right\}^{\top} \\
 &\approx [-](\hat{\theta} - \theta^*)^{\top} \left[\sum_{\alpha=1}^{\mathcal{N}} \nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} (\hat{\theta} - \theta^*)
 \end{aligned}$$

- Take the large \mathcal{N} limit

$$\lim_{\mathcal{N} \gg 1} \left[\sum_{\alpha=1}^{\mathcal{N}} \nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} \approx \mathcal{N} \left\langle \nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right\rangle_{\theta^*=\hat{\theta}} \approx [-] \mathcal{N} g(\theta^*)$$

- Thus

$$\begin{aligned}
 \lim_{\mathcal{N} \gg 1} \sum_{\alpha=1}^{\mathcal{N}} (\hat{\theta} - \theta^*)^{\top} \left[\nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x_{\alpha}) \right]_{\theta^*=\hat{\theta}} &\approx (\hat{\theta} - \theta^*)^{\top} \mathcal{N} g(\theta^*) (\hat{\theta} - \theta^*) \\
 &\approx \mathcal{N} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2
 \end{aligned}$$

Derivation

- Therefore

$$\begin{aligned}\mathcal{L}(\hat{\theta} | x) &\approx \sum_{\alpha=1}^{\mathcal{N}} \mathcal{L}(\theta^* | x_{\alpha}) + \mathcal{N} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 + [-] \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \\ &\approx \mathcal{L}(\theta^* | x) + \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2\end{aligned}$$

- Next, approximate $\overline{\mathcal{L}}(\hat{\theta} | x)$

$$\begin{aligned}\overline{\mathcal{L}}(\hat{\theta} | x) &\approx \overline{\mathcal{L}}(\theta^* | x) + (\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \overline{\mathcal{L}}(\hat{\theta} | x) \right]_{\hat{\theta}=\theta^*}^{\top} \\ &\quad + \frac{1}{2} (\hat{\theta} - \theta^*)^{\top} \left[\nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \overline{\mathcal{L}}(\hat{\theta} | x) \right]_{\hat{\theta}=\theta^*} (\hat{\theta} - \theta^*)\end{aligned}$$

- The term $\sim \mathcal{O}(\theta)$ correspond with MLE goals

$$(\hat{\theta} - \theta^*) \left[\nabla_{\hat{\theta}} \overline{\mathcal{L}}(\hat{\theta} | x) \right]_{\hat{\theta}=\theta^*}^{\top} \approx (\hat{\theta} - \theta^*)(0) \approx (0) \quad \because \theta^* = \arg \max \overline{\mathcal{L}}$$

- The term $\sim \mathcal{O}(\theta^2)$ at $\mathcal{N} \gg 1$

$$\begin{aligned}\lim_{\mathcal{N} \gg 1} \frac{1}{2} (\hat{\theta} - \theta^*)^{\top} \left[\nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \overline{\mathcal{L}}(\hat{\theta} | x) \right]_{\hat{\theta}=\theta^*} (\hat{\theta} - \theta^*) \\ \approx \frac{1}{2} (\hat{\theta} - \theta^*)^{\top} \left[\overline{\nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \mathcal{L}(\hat{\theta} | x)} \right]_{\hat{\theta}=\theta^*} (\hat{\theta} - \theta^*)\end{aligned}$$

Derivation

- continued

$$\begin{aligned} \lim_{\mathcal{N} \gg 1} \frac{1}{2} (\hat{\theta} - \theta^*)^\top \left[\nabla_{\hat{\theta}} \nabla_{\hat{\theta}} \overline{\mathcal{L}}(\hat{\theta} | x) \right]_{\hat{\theta} = \theta^*} (\hat{\theta} - \theta^*) \\ \approx \frac{1}{2} (\hat{\theta} - \theta^*)^\top ([-]g(\theta^*)) (\hat{\theta} - \theta^*) \approx [-] \frac{1}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \end{aligned}$$

- Collecting results for \mathcal{L} and $\overline{\mathcal{L}}$

$$\begin{aligned} \overline{\mathcal{L}}(\hat{\theta} | x) &\approx \overline{\mathcal{L}}(\theta^* | x) + [-] \frac{1}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \\ \mathcal{L}(\hat{\theta} | x) &\approx \mathcal{L}(\theta^* | x) + \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \end{aligned}$$

- Next, compute risks $\hat{\mathfrak{R}}$ and \mathfrak{R}

$$\begin{aligned} \mathfrak{R}(\mathcal{L}(\hat{\theta} | x)) &= [-] \mathcal{N} \mathbb{E}_{\hat{\theta}}[\overline{\mathcal{L}}] = [-] \mathcal{N} \left(\mathbb{E}_{\hat{\theta}}[\overline{\mathcal{L}}(\theta^* | x)] + [-] \frac{1}{2} \mathbb{E}_{\hat{\theta}} \left[\left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right] \right) \\ \hat{\mathfrak{R}}(\mathcal{L}(\hat{\theta} | x)) &= \sum_{\alpha=1}^{\mathcal{N}} [-] \mathcal{L}(\hat{\theta} | x_\alpha) \approx [-] \left[\mathcal{L}(\theta^* | x) + \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right] \end{aligned}$$

- Evaluate bias

$$\begin{aligned} \text{bias} &\approx \left\langle \hat{\mathfrak{R}}(\mathcal{L}(\hat{\theta} | x_\alpha)) \right\rangle_{\hat{\theta}} - \mathfrak{R}(\mathcal{L}(\hat{\theta} | x_\alpha)) \\ &\approx \left\langle [-] \left[\mathcal{L}(\theta^* | x) + \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right] \right\rangle_{\hat{\theta}} - [-] \mathcal{N} \left(\overline{\mathcal{L}}(\theta^* | x) + [-] \frac{1}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right) \end{aligned}$$

Derivation

- continued

$$\begin{aligned} \text{bias} &\approx \langle [-] \mathcal{L}(\theta^* | x) \rangle_{\hat{\theta}} + \left\langle [-] \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right\rangle_{\hat{\theta}} \\ &\quad + \mathcal{N} \mathbb{E}_{\hat{\theta}} \left[\overline{\mathcal{L}}(\theta^* | x) \right] + [-] \mathcal{N} \frac{1}{2} \mathbb{E}_{\hat{\theta}} \left[\left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right] \end{aligned}$$

- Note: 2nd term=4th term.
- 1st term is similar to 3th term

$$\begin{aligned} \langle [-] \mathcal{L}(\theta^* | x) \rangle_{\hat{\theta}} &\approx [-] \mathbb{E}_{\hat{\theta}} [\mathcal{L}(\theta^* | x)] \approx [-] \mathcal{N} \mathcal{N}^{[-]} \mathbb{E}_{\hat{\theta}} \left[\sum_{\alpha=1}^{\mathcal{N}} \mathcal{L}(\theta^* | x_{\alpha}) \right] \\ &\approx [-] \mathcal{N} \mathbb{E}_{\hat{\theta}} \left[\mathcal{N}^{[-]} \sum_{\alpha=1}^{\mathcal{N}} \mathcal{L}(\theta^* | x_{\alpha}) \right] \approx [-] \mathcal{N} \mathbb{E}_{\hat{\theta}} \left[\overline{\mathcal{L}}(\theta^* | x_{\alpha}) \right] \end{aligned}$$

- Hence

$$\begin{aligned} \text{bias} &\approx [-] \mathcal{N} \mathbb{E}_{\hat{\theta}} \left[\overline{\mathcal{L}}(\theta^* | x_{\alpha}) \right] + \left\langle [-] \frac{\mathcal{N}}{2} \left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right\rangle_{\hat{\theta}} \\ &\quad \mathcal{N} \mathbb{E}_{\hat{\theta}} \left[\overline{\mathcal{L}}(\theta^* | x_{\alpha}) \right] + [-] \mathcal{N} \frac{1}{2} \mathbb{E}_{\hat{\theta}} \left[\left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right] \\ &\approx (0) + [-] \mathcal{N} \mathbb{E}_{\hat{\theta}} \left[\left\| \hat{\theta} - \theta^* \right\|_{g(\theta^*)}^2 \right] \end{aligned}$$

Derivation

- By definition of Central Limit Theorem (CLT) and Law of Large Numbers (LLN) IID $\hat{\theta}$ distribute as Gaussians, and $(\hat{\theta} - \theta^*)^\top \mathcal{N} g(\theta^*) (\hat{\theta} - \theta^*)$ as χ^2 (prove it elsewhere?):

$$\begin{aligned}\sqrt{\mathcal{N}}(\hat{\theta} - \theta^*) &\sim \mathcal{N}(0, g^{[-]}(\theta^*)) & \mathcal{N}(\hat{\theta} - \theta^*)^\top g(\theta^*) (\hat{\theta} - \theta^*) &\sim \chi^2_{\mathcal{K}} \\ \mathbb{E}_{\hat{\theta}} \left[\mathcal{N}(\hat{\theta} - \theta^*)^\top g(\theta^*) (\hat{\theta} - \theta^*) \right] &\sim \mathcal{K}\end{aligned}$$

- Thus, bias of $\hat{\mathfrak{R}}$ is

$$\text{bias} \approx [-]_{\mathcal{K}}$$

- AIC is $\times 2$ the empirical risk (mod the bias)

$$\text{AIC} = 2(\hat{\mathfrak{R}}(\mathcal{L}(\hat{\theta} | x)) - \text{bias}) = 2\mathcal{K} - 2\mathcal{L}(\hat{\theta} | x)$$

Derivation

$$\|(\hat{\theta} - \theta^*)^\top \mathcal{N}g(\theta^*)(\hat{\theta} - \theta^*) \sim \chi^2_{\mathcal{K}}\|$$

- Useful fact: Since Fisher g is real, symmetric (and semi-definite), a rotation S can diagonalize it.
- For IID $\hat{\theta}$, Central Limit Theorem implies: $\sqrt{\mathcal{N}}(\hat{\theta} - \theta^*) \sim \mathcal{N}(0, g(\theta^*))$ (proof?)

$$\begin{aligned} G &\equiv \int d\hat{\theta} \mathcal{Z}^{[-]} e^{[-]\frac{1}{2}\|\hat{\theta} - \theta^*\|_{g(\theta^*)}^2} = \mathcal{Z}^{[-]} \int d(\hat{\theta} - \theta^*) e^{[-]\frac{1}{2}\|\hat{\theta} - \theta^*\|_{g(\theta^*)}^2} \\ &= \mathcal{Z}^{[-]} \prod_{M=1}^{\mathcal{K}} \int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2}(\hat{\theta}_I - \theta_I^*)(\hat{\theta}_J - \theta_J^*)g(\theta^*)_{IJ}} \\ &= \mathcal{Z}^{[-]} \prod_{M=1}^{\mathcal{K}} \int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2}(\hat{\theta}_I - \theta_I^*)(\hat{\theta}_J - \theta_J^*)(S^\top \text{diag}(g(\theta^*))S)_{IJ}} \\ &= \mathcal{Z}^{[-]} \prod_{M=1}^{\mathcal{K}} \int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2}(\hat{\theta}_I - \theta_I^*)(\hat{\theta}_J - \theta_J^*)(S_{IL}^\top \text{diag}(g(\theta^*))_{LN} S_{NJ})} \\ &= \mathcal{Z}^{[-]} \prod_{M=1}^{\mathcal{K}} \int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2} \sum_{I=1}^{\mathcal{K}} [S_{IJ}(\hat{\theta}_J - \theta_J^*)]^2} \\ &= \mathcal{Z}^{[-]} \prod_{M=1}^{\mathcal{K}} \left(\int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2} \text{diag}(g(\theta^*))_{MM} [S_{MN}(\hat{\theta}_N - \theta_N^*)]^2} \right) \end{aligned}$$

- Change of variable $\sqrt{\text{diag}(g(\theta^*))_{MM}} S_{MN}(\hat{\theta}_N - \theta_N^*) \rightarrow \phi_M$

$$G = \prod_{M=1}^{\mathcal{K}} \left(\int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2} \sum_{M=1}^{\mathcal{K}} \phi_M^2} \right)$$

Derivation

- Transform integration measures $(\hat{\theta}_M - \theta_M^*), \forall M \in [1, \mathcal{K}]$

$$\int d(\hat{\theta}_M - \theta_M^*) = \frac{1}{\sqrt{\text{diag}(g(\theta^*))_{MM}}} \int d(S_{NM}^\top \phi_N) = \int \det(S^\top) d\phi_N = \int d\phi_N$$

- Whence S is orthonormal (i.e. $\|S\|^2 = S^\top S = 1 = SS^\top$), and as such $\det S = 1 = \frac{1}{\det S^\top} = \det S^\top$

- Recollect $\mathcal{Z}[-] = \frac{\sqrt{\det(\mathcal{N}g(\theta^*))}}{(2\pi)^{\frac{\mathcal{K}}{2}}}$ (Jaynes¹ App E)

$$\begin{aligned} G &= \prod_{M=1}^{\mathcal{K}} \left(\int d(\hat{\theta}_M - \theta_M^*) e^{[-]\frac{1}{2} \sum_{M=1}^{\mathcal{K}} \phi_M^2} \right) \\ &= \frac{1}{\sqrt{\text{diag}(\mathcal{N}g(\theta^*))_{MM}} \mathcal{Z}} \prod_{M=1}^{\mathcal{K}} \left(\int d\phi_M e^{[-]\frac{1}{2} \sum_{M=1}^{\mathcal{K}} \phi_M^2} \right) \\ &= \frac{\sqrt{\text{diag}(\mathcal{N}g(\theta^*))_{MM}}}{\sqrt{\text{diag}(\mathcal{N}g(\theta^*))_{MM}} (2\pi)^{\frac{\mathcal{K}}{2}}} \prod_{M=1}^{\mathcal{K}} \left(\int d\phi_M e^{[-]\frac{1}{2} \sum_{M=1}^{\mathcal{K}} \phi_M^2} \right) \\ &= \frac{1}{(2\pi)^{\frac{\mathcal{K}}{2}}} \prod_{M=1}^{\mathcal{K}} \left(\int d\phi_M e^{[-]\frac{1}{2} \sum_{M=1}^{\mathcal{K}} \phi_M^2} \right) \end{aligned}$$

- Change of variable $\omega^2 := \phi_M \phi^M = \sum_{M=1}^{\mathcal{K}} \phi_M^2$ (i.e Path to polar coordinate)

$$\prod_{M=1}^{\mathcal{K}} \int d\phi_M = \prod_{M=1}^{\mathcal{K}} \int_{\mathbb{R}} d[\sqrt{\mathcal{N}}S(\hat{\theta} - \theta^*)]_M = \int_{\mathbb{R}^+} d\omega A_{\mathcal{K}-1}$$

Derivation

- Where $A_{\mathcal{K}-1} \equiv$ Surface Area of $(\mathcal{K} - 1)$ -Shell $= \frac{2\pi^{\mathcal{K}/2}}{\Gamma(\mathcal{K}/2)} \omega^{\mathcal{K}-1}$

$$G = \frac{1}{(2\pi)^{\frac{\mathcal{K}}{2}}} \int_{\mathbb{R}^+} d\omega A_{\mathcal{K}-1} e^{[-]\frac{1}{2}\omega^2} = \frac{1}{(2\pi)^{\frac{\mathcal{K}}{2}}} \int_{\mathbb{R}^+} d\omega \frac{2\pi^{\mathcal{K}/2}}{\Gamma(\mathcal{K}/2)} \omega^{\mathcal{K}-1} e^{[-]\frac{1}{2}\omega^2}$$

- Change variable $\omega^2 := \zeta$

$$\begin{aligned} G &:= \frac{1}{(2\pi)^{\frac{\mathcal{K}}{2}}} \int_{\mathbb{R}^+} d\sqrt{\zeta} \frac{2\pi^{\mathcal{K}/2}}{\Gamma(\mathcal{K}/2)} \zeta^{\frac{1}{2}(\mathcal{K}-1)} e^{[-]\frac{1}{2}\zeta} \\ &= \frac{2\pi^{\mathcal{K}/2}}{\Gamma(\mathcal{K}/2)(2\pi)^{\frac{\mathcal{K}}{2}}} \int_{\mathbb{R}^+} \frac{d\zeta}{2\sqrt{\zeta}} \zeta^{\frac{1}{2}(\mathcal{K}-1)} e^{[-]\frac{1}{2}\zeta} \\ &= \frac{1}{\Gamma(\mathcal{K}/2)(2)^{\frac{\mathcal{K}}{2}}} \int_{\mathbb{R}^+} d\zeta \zeta^{\frac{\mathcal{K}}{2}-1} e^{[-]\frac{1}{2}\zeta} \\ &= \int_{\mathbb{R}^+} d\zeta \frac{\zeta^{\frac{\mathcal{K}}{2}-1} e^{[-]\frac{1}{2}\zeta}}{\Gamma(\mathcal{K}/2)(2)^{\frac{\mathcal{K}}{2}}} \\ &\equiv \int d\zeta \chi^2[\mathcal{K}] \end{aligned}$$

- Thus, chi-squared of degree \mathcal{K} is

$$\therefore \chi^2[\mathcal{K}] \equiv \frac{\zeta^{\frac{\mathcal{K}}{2}-1} e^{[-]\frac{1}{2}\zeta}}{\Gamma(\mathcal{K}/2)(2)^{\frac{\mathcal{K}}{2}}}$$

Derivation

$$\mathbb{E}_{\hat{\theta}} \left[(\hat{\theta} - \theta^*)^\top \mathcal{N}g(\theta^*)(\hat{\theta} - \theta^*) \right] \sim \mathcal{K}$$

- Introduce χ^2 of degree 1

$$\chi^2[\mathcal{K} = 1] \equiv \frac{1}{\sqrt{2\pi\omega e^\omega}}$$

- Compute the moment generator \mathfrak{m} for a Γ distribution, on $\omega \in [-\infty, \infty]$ (not $\text{dom}(\chi^2) = (0, \infty)$)

$$\mathfrak{m} \equiv \int_{\mathbb{R}} d\omega e^{t\omega} \chi^2(\omega) = \int_{\mathbb{R}} d\omega e^{t\omega} \frac{1}{\sqrt{2\pi\omega e^\omega}}$$

- Change of variable $\omega \mapsto \tilde{\omega}^2$

$$\begin{aligned} \mathfrak{m} &= 2 \cdot \frac{1}{\sqrt{2\pi}} \int_0^\infty d\tilde{\omega}^2 e^{t\tilde{\omega}^2} \frac{1}{\sqrt{\tilde{\omega}^2 e^{\tilde{\omega}^2}}} = 2 \cdot \frac{1}{\sqrt{2\pi}} \int_0^\infty d\tilde{\omega} e^{(-)\frac{1}{2}(1-2t)\tilde{\omega}^2} \\ &= 2 \cdot \frac{1}{\sqrt{2\pi}} \frac{1}{2} \sqrt{\frac{\pi}{\frac{1}{2}(1-2t)}} = \sqrt{\frac{\pi}{(1-2t)}} \end{aligned}$$

Derivation cont

- For a random IID $\rho \equiv \sum_{I=1}^{\mathcal{K}} \omega_I$ such that $\rho \sim p(\rho) \equiv \prod_{I=1}^{\mathcal{K}} p(\omega_I)$

$$\begin{aligned} \mathfrak{m}_{\mathcal{K}} &= \int d\rho_{\mathbb{R}^+} e^{t\rho} p(\rho) = \int d\omega_1 \dots \int d\omega_{\mathcal{K}} e^{t \sum_{J=1}^{\mathcal{K}} \omega_J} \prod_{I=1}^{\mathcal{K}} p(\omega_I) \\ &= \prod_{J=1}^{\mathcal{K}} \left[\int d\omega_J e^{t\omega_J} p(\omega_J) \right] = \left[\int d\omega_1 e^{t\omega_1} p(\omega_{J=1}) \right]^{\mathcal{K}} = \mathfrak{m}^{\mathcal{K}} \end{aligned}$$

- Apply on χ^2

$$\mathfrak{m}_{\mathcal{K}} = \left[\int d\omega_1 e^{t\omega_1} \chi^2 \right]^{\mathcal{K}} = \left[(1 - 2t) \right]^{(-) \frac{\mathcal{K}}{2}}$$

- Useful fact

$$\begin{aligned} \left[\frac{\partial}{\partial \xi} \mathfrak{m}_{\xi}[\chi^2(X)] \right]_{\xi=0} &= \left[\frac{\partial}{\partial \xi} \int dx e^{\xi x} \chi^2 \right]_{\xi=0} = \left[\int dx \frac{\partial}{\partial \xi} e^{\xi x} \chi^2 \right]_{\xi=0} = \left[\int dx x e^{\xi x} \chi^2 \right]_{\xi=0} \\ &= \mathbb{E}_{x \sim \chi^2}(X) \end{aligned}$$

- Let $x \equiv \omega = \|\hat{\theta} - \theta^*\|_{\mathcal{N}_{g(\theta^*)}}^2$

$$\begin{aligned} \mathbb{E}_{x \sim \chi^2}(X) &= \left\langle \|\hat{\theta} - \theta^*\|_{\mathcal{N}_{g(\theta^*)}}^2 \right\rangle = \left[\frac{\partial}{\partial \xi} \mathfrak{m}_{\xi}[\chi^2(x = \|\hat{\theta} - \theta^*\|_{\mathcal{N}_{g(\theta^*)}}^2)] \right]_{\xi=0} \\ &= \left[\frac{\partial}{\partial \xi} \mathfrak{m}_{\xi}[\chi^2(x = \|\hat{\theta} - \theta^*\|_{\mathcal{N}_{g(\theta^*)}}^2)] \right]_{\xi=0} = \left[\frac{\partial}{\partial \xi} (1 - 2t)^{(-) \frac{1}{2} \mathcal{K}} \right]_{\xi=0} \\ &= \left[\mathcal{K} (1 - 2\xi)^{(-) \frac{3}{2} \mathcal{K}} \right]_{\xi=0} = \mathcal{K} \end{aligned}$$

Derivation

$$\text{"BIC} = \mathcal{K} \log \mathcal{N} - \mathcal{L}(\hat{\theta} | x) \text{"}$$

- Assume models $m_J, \forall J \in [1, \mathcal{K}]$ *Model Probability* defines to be

$$\begin{aligned} p(m_J | x) &= \frac{p(m_J, x)}{p(x)} = \int d\theta_J \frac{p(m_J, \theta_J, x)}{p(x)} = \int d\theta_J \frac{p(x | m_J, \theta_J) p(\theta_J | m_J) p(m_J)}{p(x)} \\ &= \frac{p(m_J)}{p(x)} \int d\theta_J p(x | m_J, \theta_J) p(\theta_J | m_J) = \frac{p(m_J)}{p(x)} \int d\theta_J \mathcal{L}(\hat{\theta}) \end{aligned}$$

- Let $p(m_J | x)$ assumes the role of likelihood. Energy E defines to be

$$E = [-] \log p(m_J | x) = [-] \log \frac{p(m_J)}{p(x)} \int d\theta_J \mathcal{L}(\hat{\theta}) = [-] \log \frac{p(m_J)}{p(x)} + [-] \log \int d\theta_J \mathcal{L}(\hat{\theta})$$

- Log-Likelihood *completed* with $\{\theta_J\}$ defines to be

$$\mathcal{L}(\theta_J) \equiv \log p(x | m_J \theta_J)$$

- Taylor expand around $\sim \hat{\theta}_J$

$$\begin{aligned} \mathcal{L}(\theta_J) &\approx \mathcal{L}(\hat{\theta}_J) + (\theta_J - \hat{\theta}_J)^\top \left(\nabla_{\theta_J} \mathcal{L}(\theta_J) \right)_{\theta=\hat{\theta}} + \frac{1}{2} (\theta_J - \hat{\theta}_J)^\top \left(\nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J) \right)_{\theta=\hat{\theta}} (\theta_J - \hat{\theta}_J) \\ &\approx \mathcal{L}(\hat{\theta}_J) + (\theta_J - \hat{\theta}_J)^\top (0) + \frac{1}{2} (\theta_J - \hat{\theta}_J)^\top \left(\nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J) \right)_{\theta=\hat{\theta}} (\theta_J - \hat{\theta}_J) \end{aligned}$$

- Where $\nabla_{\theta_J} \mathcal{L}(\theta_J) = 0$ for $\hat{\theta}_J \equiv \arg \max_{\theta} \mathcal{L}$

Derivation cont

- Term $\sim \mathcal{O}(\theta^2)$

$$\begin{aligned}\left(\nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J)\right)_{\theta=\hat{\theta}} &= \left[\frac{\mathcal{N}}{\mathcal{N}} \nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J)\right]_{\theta=\hat{\theta}} = \left[\frac{\mathcal{N}}{\mathcal{N}} \sum_{\alpha=1}^{\mathcal{N}} \nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J)\right]_{\theta=\hat{\theta}} \\&= \mathcal{N} \left[\frac{1}{\mathcal{N}} \sum_{\alpha=1}^{\mathcal{N}} \nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J)\right] \approx \mathcal{N} \left[\langle \nabla_{\theta_J} \nabla_{\theta_J} \mathcal{L}(\theta_J) \rangle_{x \sim p(x\theta)}\right]_{\theta=\hat{\theta}} \\&\approx [-] \mathcal{N} \left[g(\theta)\right]_{\theta=\hat{\theta}} \approx [-] \mathcal{N} g(\hat{\theta})\end{aligned}$$

- Thus, log-likelihood around $\hat{\theta}_J$ approximates to be

$$\mathcal{L}(\theta_J) \approx \mathcal{L}(\hat{\theta}_J) + \frac{\mathcal{N}}{2} (\theta - \hat{\theta})^\top g(\hat{\theta}) (\theta - \hat{\theta}) \approx \mathcal{L}(\hat{\theta}_J) + [-] \frac{1}{2} \|\theta - \hat{\theta}\|_{g(\hat{\theta})}^2$$

- Exponentiate that approximation \mathcal{L} and put it into E

$$\begin{aligned}[-] \log p(m_J | x) &\approx [-] \log \frac{p(m_J)}{p(x)} + [-] \log \int d\theta_J e^{\mathcal{L}(\hat{\theta}_J) + [-] \frac{1}{2} \|\theta - \hat{\theta}\|_{g(\hat{\theta})}^2} \\&\approx [-] \log \frac{p(m_J)}{p(x)} + [-] \log \left[e^{\mathcal{L}(\hat{\theta}_J)} \int d\theta_J e^{[-] \frac{1}{2} \|\theta - \hat{\theta}\|_{g(\hat{\theta})}^2} \right] \\&\approx [-] \log \frac{p(m_J)}{p(x)} + [-] \mathcal{L}(\hat{\theta}_J) + [-] \log \int d\theta_J e^{[-] \frac{1}{2} \|\theta - \hat{\theta}\|_{g(\hat{\theta})}^2}\end{aligned}$$

Derivation cont

- For Multivariate Gaussian

$$\int d\theta_J e^{[-]\frac{1}{2}\|\theta-\hat{\theta}\|_{g(\hat{\theta})}^2} = \frac{1}{\sqrt{2\pi\det(\mathcal{N}g(\hat{\theta}))}\mathcal{K}} = \frac{1}{\sqrt{2\pi\mathcal{N}\det(g(\hat{\theta}))}\mathcal{K}}$$

- Energy E approximates to be

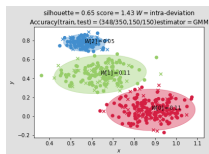
$$\begin{aligned} E &\approx [-]\log \frac{p(m_J)}{p(x)} + [-]\mathcal{L}(\hat{\theta}_J) + [-]\log \frac{1}{\sqrt{2\pi\mathcal{N}\det(g(\hat{\theta}))}\mathcal{K}} \\ &\approx [-]\log \frac{p(m_J)}{p(x)} + [-]\mathcal{L}(\hat{\theta}_J) + [-]\log \frac{1}{\sqrt{2\pi\det(g(\hat{\theta}))}\mathcal{K}} + [-]\log \frac{1}{\sqrt{\mathcal{N}}\mathcal{K}} \end{aligned}$$

- Keep terms $\sim \mathcal{O}(\mathcal{N})$. Disregard the rest insensitive, or subleading, to batch size \mathcal{N}

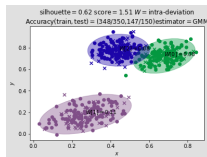
$$\text{BIC} = [-]\mathcal{L}(\hat{\theta}_J) + [-]\log \frac{1}{\sqrt{\mathcal{N}}\mathcal{K}} = \mathcal{K} \log \mathcal{N} + [-]\mathcal{L}(\hat{\theta}_J)$$

Case Study: GMM

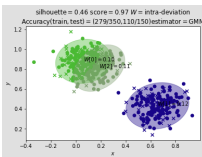
- 4 types of Gaussian clusters estimated by GMM² (with #labels = fixed.)



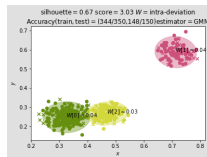
(a) Sizes



(b) Slender

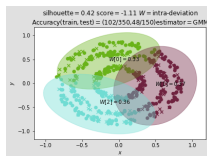


(c) Spherical

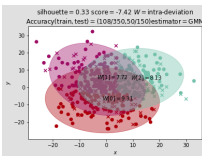


(d) Populations³

- 2 cases where GMM fails to cluster correctly.



(e) Rings



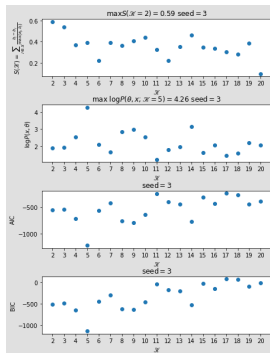
(f) Blur

²<https://github.com/scikit-learn/scikit-learn/tree/main/sklearn>

³#points for each sub-cluster C_I are distributed by random weights

Case Study: GMM continued

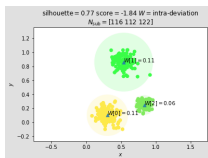
- Evaluating different scores on a parameter set of models $\mathcal{K} \equiv \# \text{initial clusters}$.



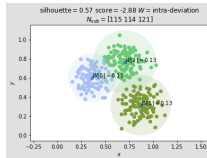
- seed $\equiv \# \text{clusters used for data generation}$

Case Study: K-Means

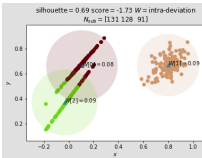
- K-means cluster different geometries of Gaussian blobs.



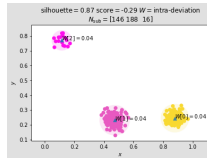
(g) Sizes



(h) Identical

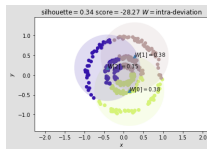


(i) Slender

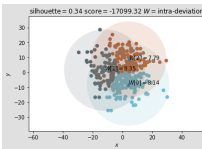


(j) Populations

- K-means failed to cluster rings (Gaussian error bars) and blurry blobs.



(k) Rings



(l) Blur