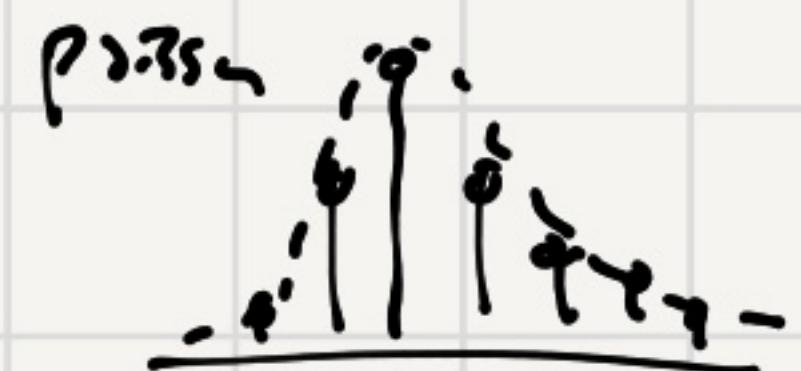
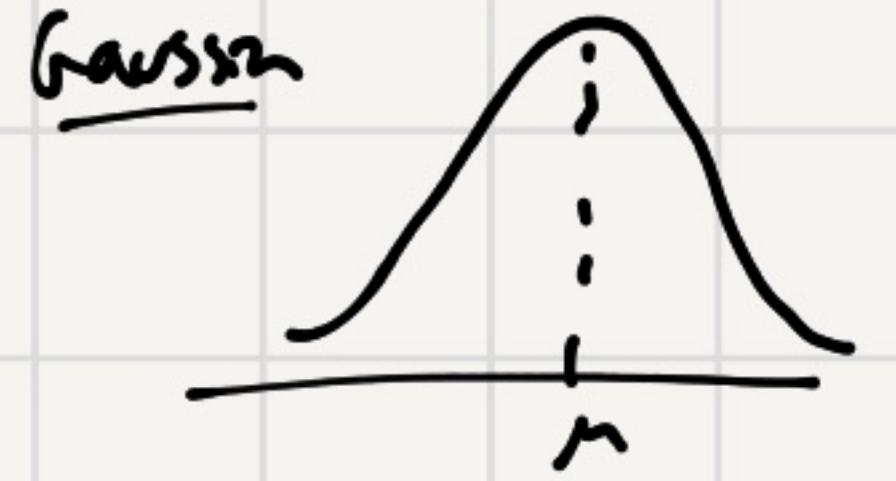
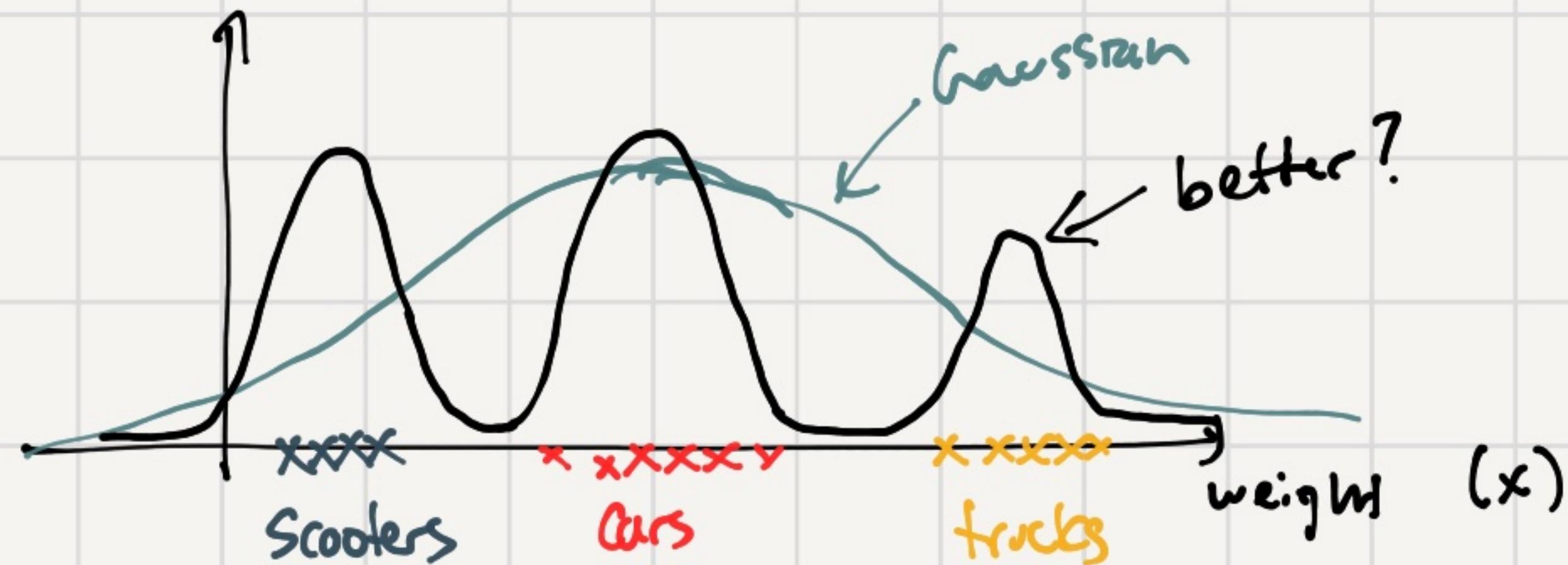


Lecture 4: Mixture Models & Clustering.



So far, our probability models have one mode (peak).

e.g. Bridge sensor - measures weight of a vehicle.



If the data is more complicated, then a Gaussian can't tell the whole story.

Gaussian Mixture Model (GMM)

"hidden state" Z - we don't directly observe this

Two r.v.s = {
 1) type of vehicle: $Z \in \{ \text{Scooter, car, truck} \}$
 (1), (2), (3)
 2) "observation" X
 weight of vehicle, conditioned on vehicle type.
 $p(x|z=j) = N(x|\mu_j, \sigma_j^2)$
 each vehicle type has its own Gaussian (mean, variance).}

Generative process:

- 1) Sample Z (type of vehicle)
- 2) Sample $X|Z$ (weight of specific type)

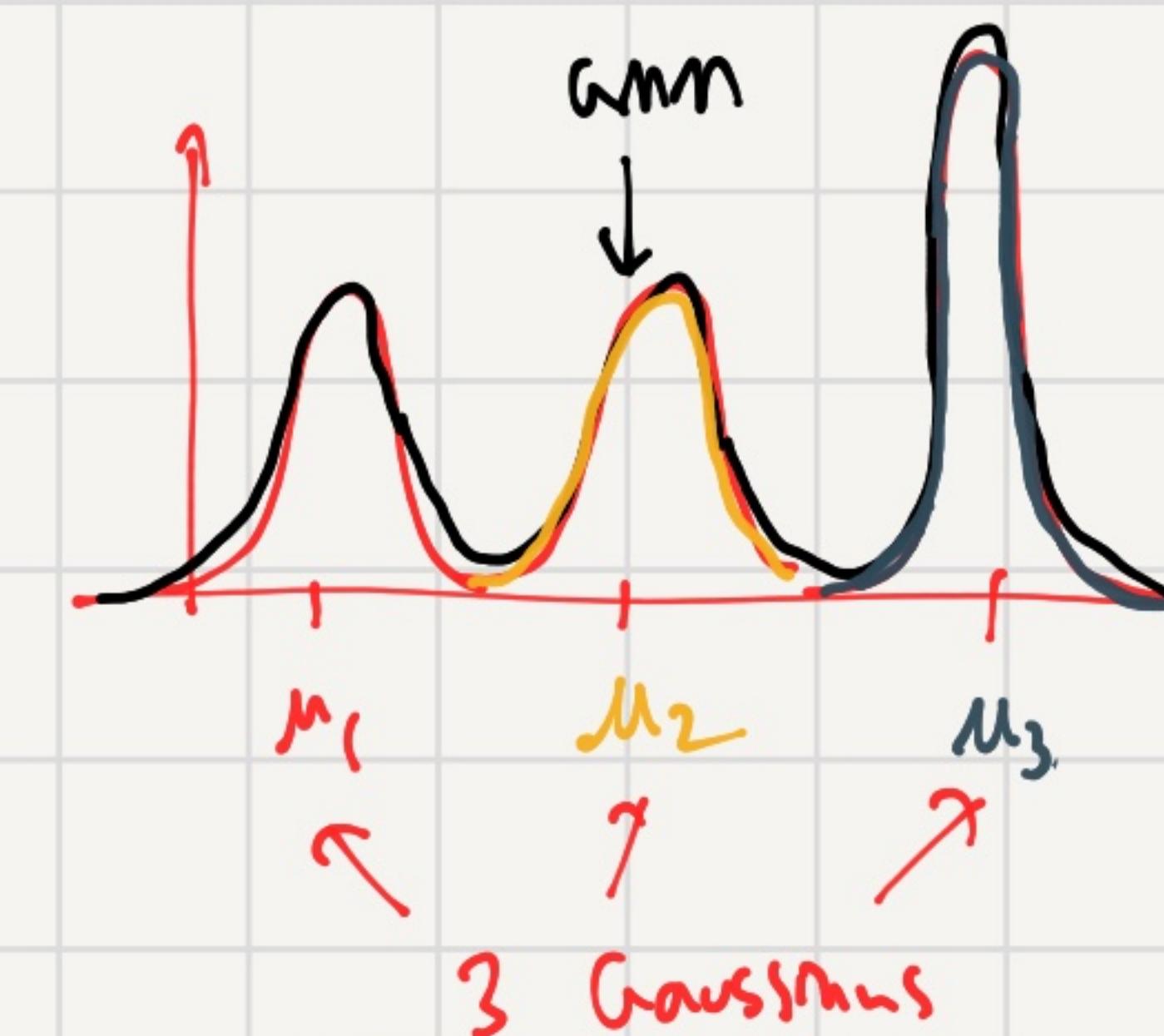
Note: we only observe X (weight). We don't see Z .

PDF of X

$$p(x) = \sum_z p(x, z) = \sum_{j=1}^K p(x|z=j) p(z=j)$$

$$\Rightarrow p(x) = \sum_{j=1}^K \pi_j N(x|\mu_j, \sigma_j^2) \quad \text{GMM}$$

Component weight mixture component (one peak)



"Weighted sum of Gaussian pdfs"

(note: not the same as a weighted sum of Gaussian r.v.s)

Clustering w/ GMMs

Dataset $D = \{x_1, \dots, x_n\}$. Assume K clusters

estimate a GMM from D :

- 1) Gaussian component \rightarrow location μ_j , extent σ_j^2
- 2) component weight \rightarrow probability of cluster π_j (# of samples)
- 3) Cluster assignment of each sample x_i (z_i)

Axtoni's Hack

let: $z_i \in \{1, \dots, K\}$ = assignment of x_i to a cluster.

- treat z_i 's as a parameter (optimize them)

- objective: maximize the joint LL of x, z

$$\hat{\theta} = \underset{\theta, z}{\operatorname{argmax}} \sum_{i=1}^n \log p(x_i, z_i) = \underset{\theta, z}{\operatorname{argmax}} \sum_i \log p(x_i | z_i) + \log p(z_i)$$

$$= \underset{\theta, z}{\operatorname{argmax}} \sum_i \log N(x_i | \mu_{z_i}, \sigma_{z_i}^2) + \log \pi_{z_i}$$

Indicator variable trick

let $z_{ij} = \begin{cases} 1, & z_i=j \text{ (} x_i \text{ is assigned to cluster } j \text{)} \\ 0, & \text{otherwise} \end{cases}$

$$p(z_i) = \prod_{j=1}^K \pi_j^{z_{ij}} \Rightarrow \log p(z_i) = \sum_{j=1}^K z_{ij} \log \pi_j$$

when
Select $\log \pi_j \wedge z_{ij}=1$

Selects π_j
when $z_{ij}=1$

$$p(x_i | z_i) = \prod_{j=1}^K N(x_i | \mu_j, \sigma_j^2)^{z_{ij}} \Rightarrow \log p(x_i | z_i) = \sum_{j=1}^K z_{ij} \log N(x_i | \mu_j, \sigma_j^2)$$

Select the j^{th}
Gaussian

$$\Rightarrow \underset{\theta, z}{\operatorname{argmax}} \sum_{i=1}^n \sum_{j=1}^K z_{ij} (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2))$$

objective J.
 $\{\mu_j, \sigma_j^2, \pi_j\}$ z_i

- variables depend on each other: alternating maximization

- 1) optimize z with θ fixed.

- 2) optimize θ with z fixed.

1) Find z_i w/ $\theta = \{\mu_i, \sigma_j^2, \pi_i\}$ fixed.

$$J = \sum_{i=1}^n \sum_{j=1}^K z_{ij} (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2))$$

each z_i is independent of other z_i 's:

$$\underset{z_{ij}}{\operatorname{argmax}} \sum_{j=1}^K z_{ij} (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2))$$

Recall $z_{ij} \in \{0, 1\}$, only one $z_{ij}=1$ for a given i .

\Rightarrow since only one term in the sum can be selected w/ $z_{ij}=1$, then just pick the max term.

$$z_i = \underset{j}{\operatorname{argmax}} \log \pi_j + \log N(x_i | \mu_j, \sigma_j^2)$$

Select j w/ largest joint log-likelihood

2) Find $\{\mu, \sigma^2, \pi\}$ given z fixed.

$$\hat{\mu}_j = \underset{\mu_j}{\operatorname{argmax}} \sum_{i=1}^n \sum_{j=1}^K z_{ij} (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2)) - \frac{1}{2\sigma_j^2} (x_i - \mu_j)^2 + \dots$$

$$= \underset{\mu_j}{\operatorname{argmax}} \sum_{i=1}^n z_{ij} \left(-\frac{1}{2\sigma_j^2} (x_i - \mu_j)^2 \right)$$

$$\frac{\partial}{\partial \mu_j} = \sum_i z_{ij} \left(-\frac{1}{2\sigma_j^2} 2(x_i - \mu_j) (-1) \right) = 0$$

$$\sum_i z_{ij} x_i - \left[\sum_i z_{ij} \right] \mu_j = 0$$

$$\Rightarrow \hat{\mu}_j = \frac{1}{\sum_{i=1}^n z_{ij}} \sum_{i=1}^n z_{ij} x_i$$

mean of samples assigned to j.
sum of samples assigned to j.
samples assigned to j.

Similarly

$$\Rightarrow \hat{\sigma}_j^2 = \frac{\sum_{i=1}^N z_{ij} (x_i - \hat{\mu}_j)^2}{\sum_{i=1}^N z_{ij}}$$

] variance of samples assigned to j

$$\Rightarrow \hat{\pi}_j = \frac{1}{N} \sum_{i=1}^N z_{ij}$$

] fraction of samples assigned to j.

3) iterate (1) & (2) until convergence.

Notes: • this 2-step iterations always maximizes our original objective $J \rightarrow$ converge (to a local max)

• if $\hat{\sigma}_j^2 = c$, $\hat{\pi}_j = \frac{1}{K}$, then this is the K-means algorithm (Lloyd's algorithm)

$$1) z_{ij} = \underset{j}{\operatorname{argmin}} (x_i - \mu_j)^2$$

$$2) \mu_j = \frac{1}{\sum z_{ij}} \sum_i z_{ij} x_i$$

• need some initial values of z or θ . (solution depends on initialization)

• Not maximizing the actual data log-likelihood

$$\log p(D) = \sum_{i=1}^N (\log p(x_i)).$$

Expectation-Maximization (EM) Algorithm (Dempster, Laird, Rubin) 1977
73,000 citations

MLE when there are hidden variables

X = observation r.v.

Z = hidden r.v.

$$\text{Joint LL: } p(X, Z) = p(X|Z)p(Z). \text{ Data LL: } p(X) = \sum_Z p(X|Z)p(Z)$$

[Goal]

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(X) = \underset{\theta}{\operatorname{argmax}} \log \sum_Z p(X|Z)p(Z)$$

Key observation: if we have both X, Z, then it's easy.

$$\operatorname{argmax} \log p(X, Z) = \operatorname{argmax} \log p(X|Z) + \log p(Z)$$

(same as step 2 in my algorithm)

- guess the values of Z in probabilistic way.
 - 1) use expected value of Z given our current model & data $\Rightarrow \hat{Z}$
 - 2) maximize $\log p(X, \hat{Z})$ to get a new model.
 - 3) repeat.

Formally

- 0) Select initial model $\hat{\theta}^{(\text{old})}$
- 1) E-step: $Q(\theta; \hat{\theta}^{(\text{old})}) = \mathbb{E}_{Z|X, \hat{\theta}^{(\text{old})}} [\log p(X, Z | \theta)]$

↑ Q-function param ↑ current model (fixed) joint LL using param θ .
conditional expectation using current model.
- 2) M-step: $\hat{\theta}^{(\text{new})} = \underset{\theta}{\operatorname{argmax}} Q(\theta; \hat{\theta}^{(\text{old})})$
- 3) $\hat{\theta}^{(\text{old})} \leftarrow \hat{\theta}^{(\text{new})}$, iterate (1) & (2)

EM for GMMs

Joint LL:

$$\log p(x, z | \theta) = \sum_{i=1}^N \sum_{j=1}^K z_{ij} (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2))$$

1) E-step:

$$\begin{aligned} Q(\theta; \hat{\theta}^{(old)}) &= \mathbb{E}_{z|x, \hat{\theta}^{(old)}} [\log p(x, z | \theta)] \\ &= \sum_{i=1}^N \sum_{j=1}^K \mathbb{E}_{z|x, \hat{\theta}^{(old)}} [z_{ij}] (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2)) \\ &= \sum_i \sum_j \hat{z}_{ij} (\log \pi_j + \log N(x_i | \mu_j, \sigma_j^2)) \end{aligned}$$

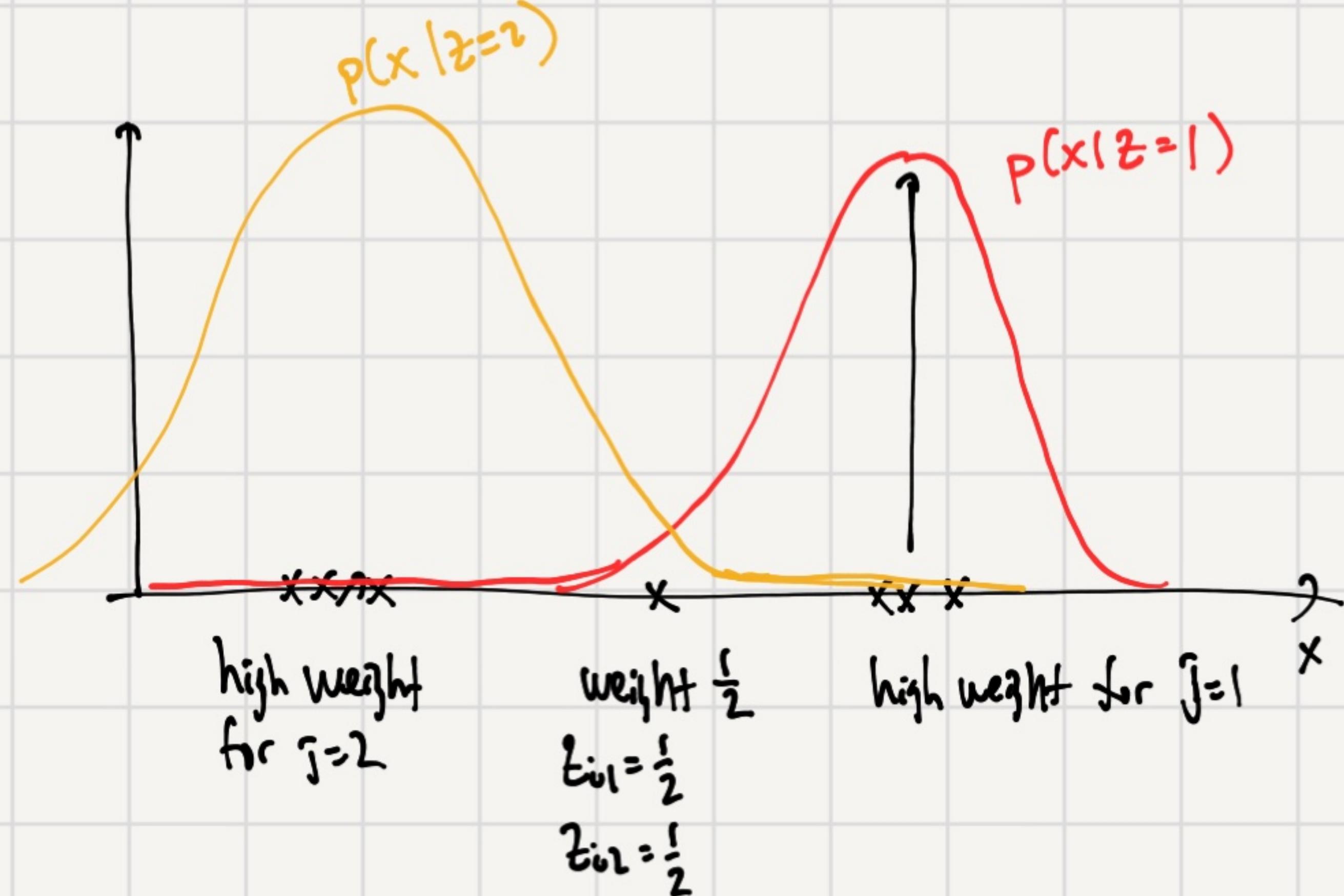
Similar to joint LL, but with \hat{z}_{ij} instead of z_{ij}

$$\begin{aligned} \hat{z}_{ij} &= \mathbb{E}_{z|x, \hat{\theta}^{(old)}} [z_{ij}] \quad \downarrow \text{expectation of an indicator variable (PSI-5)} \\ &= p(z_{ij}=1 | x, \hat{\theta}^{(old)}) \\ &= \frac{p(x | z_{ij}=1) p(z_{ij}=1)}{p(x)} \quad \downarrow \text{Bayes' Rule} \\ &= \frac{p(x_i | z_{ij}=1) p(x_{\setminus i}) p(z_{ij}=1)}{p(x_i) p(x_{\setminus i})} \quad \downarrow x_i's are independent \\ &\quad p(x) = p(x_i) p(x_{\setminus i}) \end{aligned}$$

$$\begin{aligned} \hat{z}_{ij} &= \frac{\hat{\pi}_j N(x_i | \hat{\mu}_j, \hat{\sigma}_j^2)}{\sum_k \hat{\pi}_k N(x_i | \hat{\mu}_k, \hat{\sigma}_k^2)} \quad \leftarrow \text{calculated with } \hat{\theta}^{(old)} \\ &= \sum \hat{\pi}_j, \hat{\sigma}_j^2, \hat{\mu}_j \} \quad \leftarrow \text{"soft assignment" of } x_i \text{ to cluster } j. \\ \hat{z}_{ij} &= p(z_{ij}=1 | x_i, \hat{\theta}^{(old)}) \quad \leftarrow \text{posterior prob } [0, 1] \end{aligned}$$

2) M-step: same as step 2 but with $z_{ij} \rightarrow \hat{z}_{ij}$:

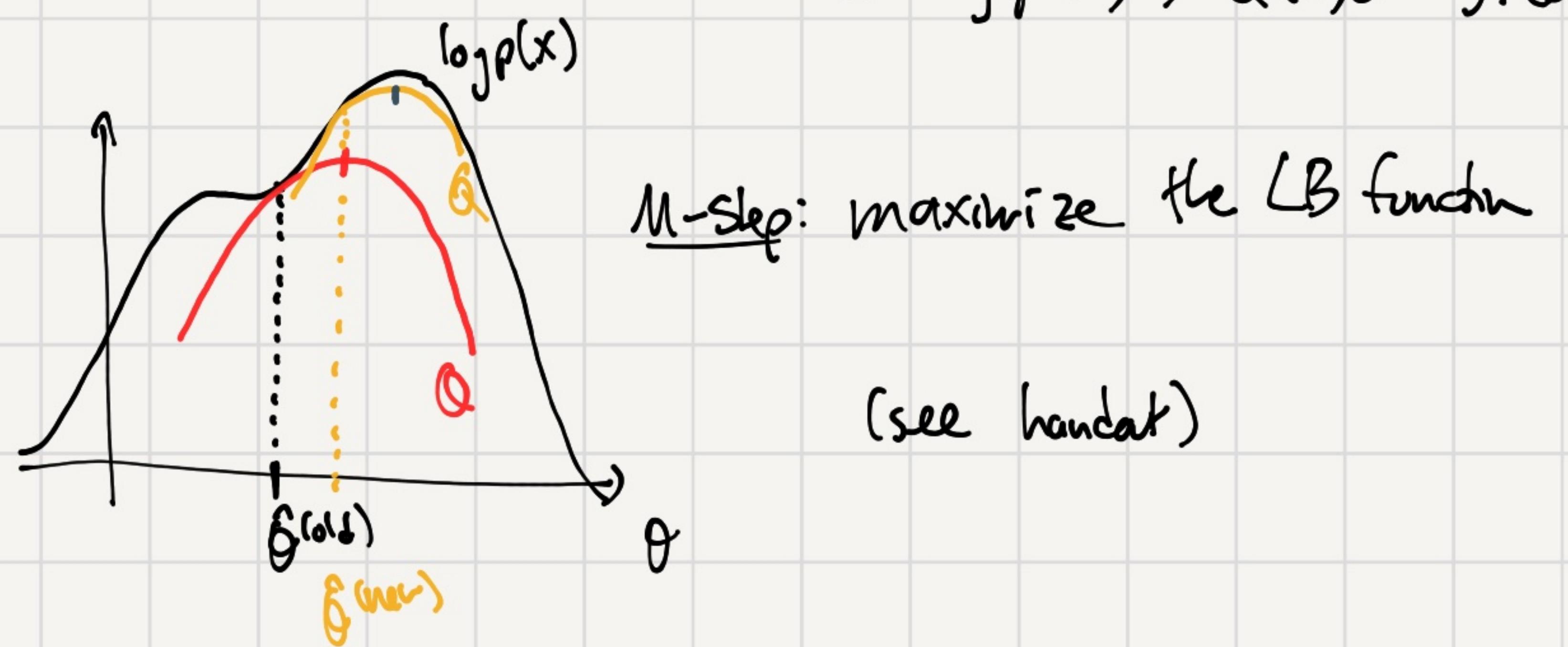
$$\left\{ \begin{array}{l} \hat{\mu}_j = \frac{1}{N_j} \sum_{i=1}^N \hat{z}_{ij} x_i \quad \leftarrow \text{Sample mean w/ weights based on soft assignments.} \\ N_j = \sum_{i=1}^N \hat{z}_{ij} \quad \leftarrow \text{total weight of samples assigned to } j. \\ \hat{\sigma}_j^2 = \frac{1}{N_j} \sum_{i=1}^N \hat{z}_{ij} (x_i - \hat{\mu}_j)^2 \quad \leftarrow \text{same} \\ \hat{\pi}_j = \frac{N_j}{N} \quad \leftarrow \text{fraction of weight assigned to cluster } j. \end{array} \right.$$



Notes:

- 1) General: EM is a general framework for MLE or models w/ hidden variables
- 2) Convergence: each iteration increases data LL ($p(x)$)
→ converges to a local max.
(but it could be slow)
- 3) Initialization: different init → different $\hat{\theta}$
pick $\hat{\theta}$ w/ largest data LL.
- 4) Interpretation:

E-step: construct a lower bound function
 $\log p(x) \geq Q(\theta; \hat{\theta}^{(t)}) + \text{const}$



M-step: maximize the LB function

(see handout)