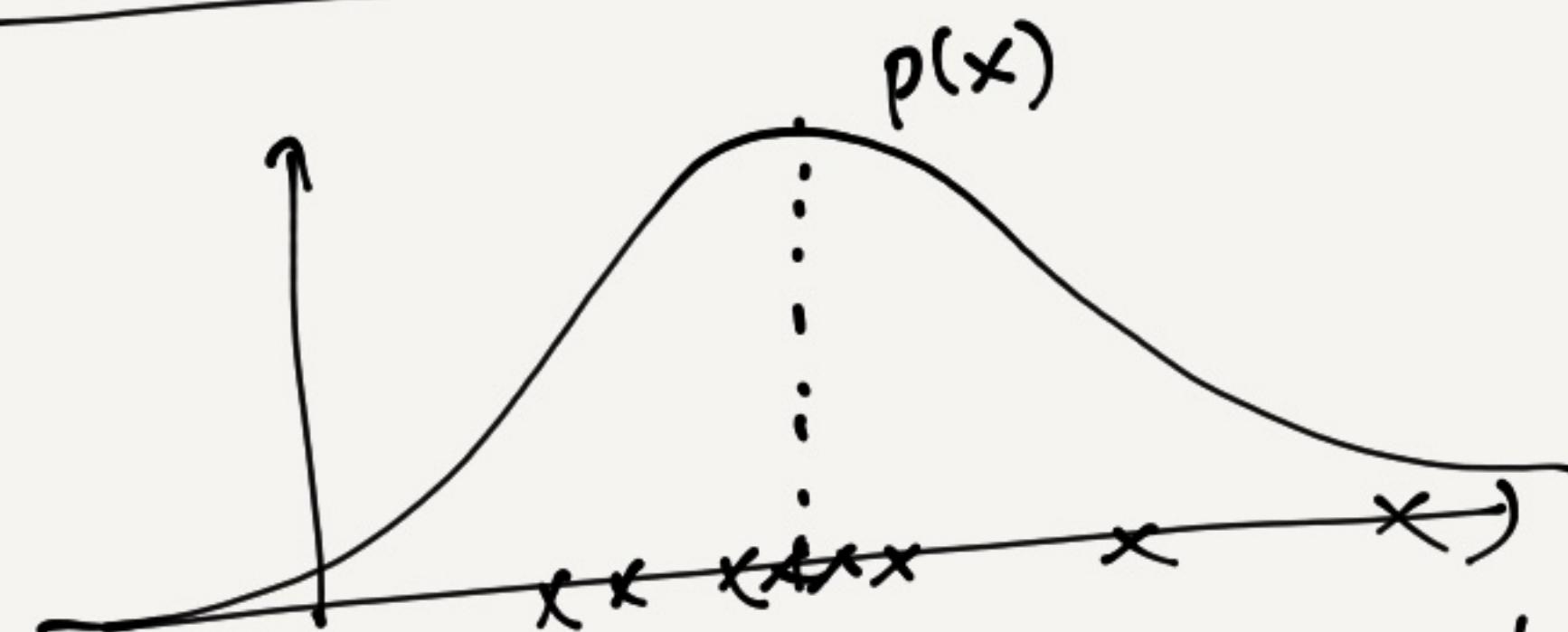


# Lecture 4

## Mixture Models & Clustering

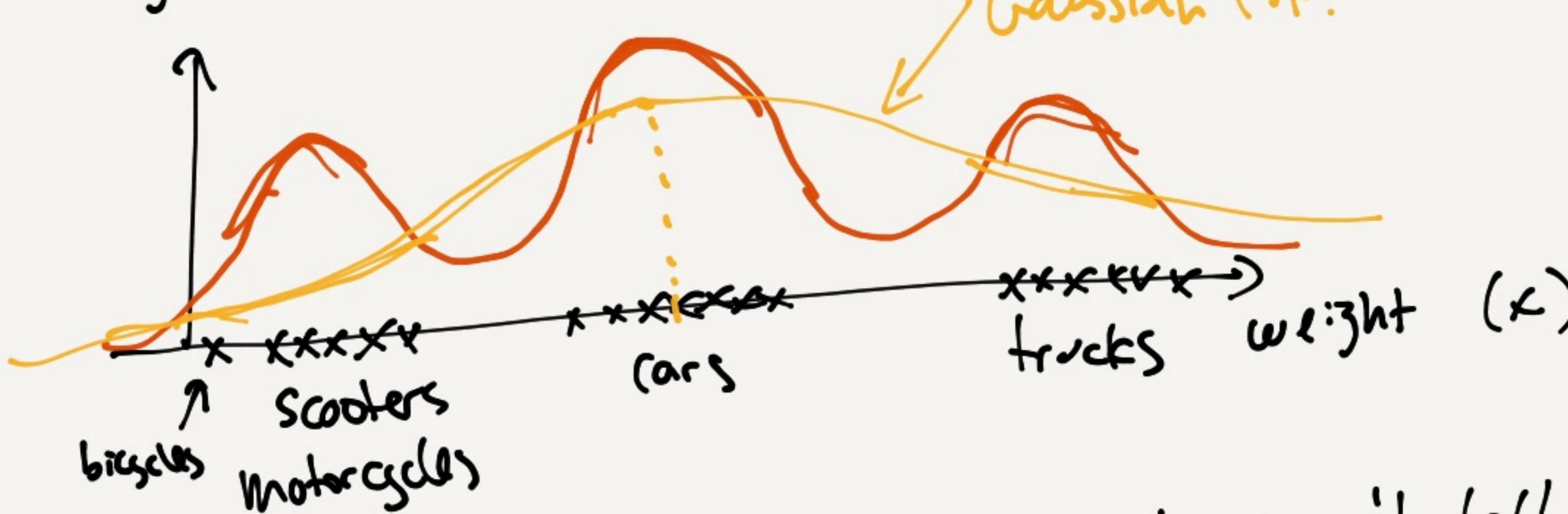


So far our probabilistic models have one mode (peak)

$$\text{e.g. } p(x) = N(x|\mu, \sigma^2)$$

What if the data is more complicated?

e.g. Bridge sensor - measure weight of vehicles



Gaussian doesn't fit well - it doesn't tell the whole story.

## Gaussian Mixture Model (GMM)

two random variables:

i)  $z$  = hidden state: e.g. vehicle type  
 $\nwarrow$  not observed  
 $z \in \{ \text{scooter, car, truck} \}$

$$p(z=j) = \pi_j, \sum_{j=1}^K \pi_j = 1$$

$\nwarrow$  prior probability of state  $z=j$  occurring.

ii)  $x$  = observation e.g. weight measured by sensor  
likelihood of  $x$  is conditioned on the current state  $z$ .

$$p(x|z=j) = N(x|\mu_j, \sigma_j^2)$$

$\nwarrow$  a Gaussian for each vehicle type.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Generative process: 1) sample a  $z$  according to  $p(z)$

2) sample  $x$  given  $z$ ,  $p(x|z)$

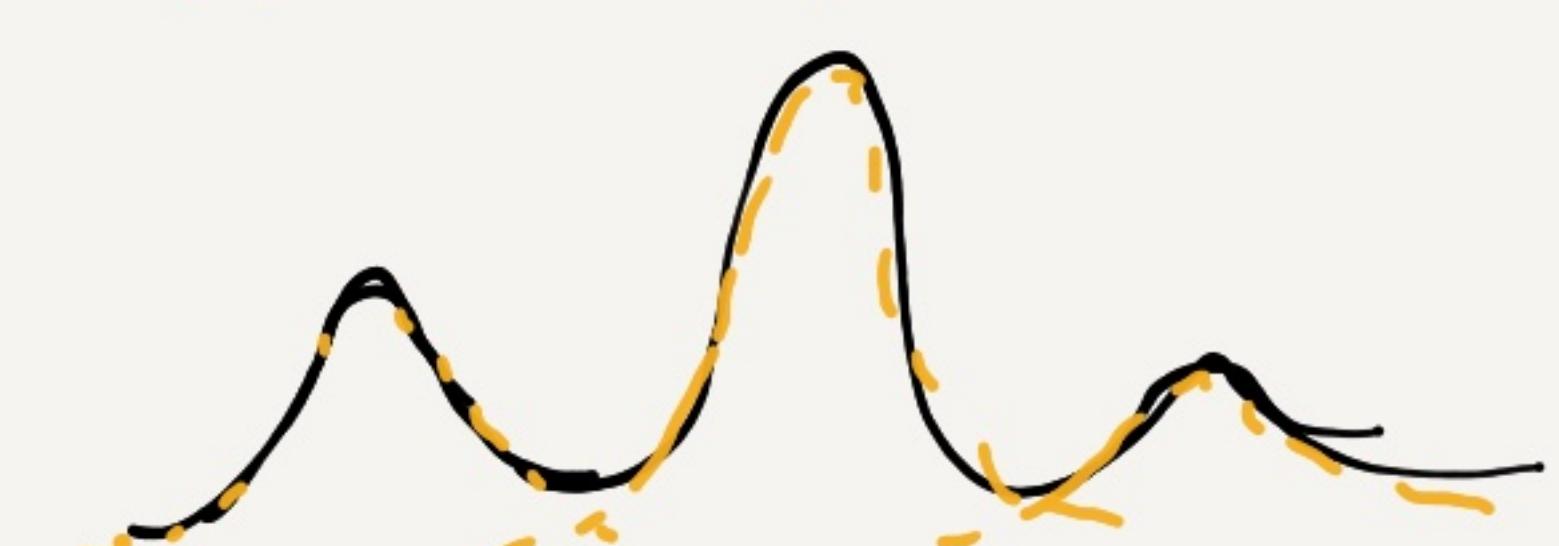
Note: we never see  $z$ ! only see  $x$ !

Likelihood of  $x$ :

$$p(x) = \sum_j p(x, z=j) = \sum_j p(x|z=j) p(z=j)$$

$$p(x) = \sum_j \pi_j N(x|\mu_j, \sigma_j^2)$$

component weight (prior)  
mixture component



## Clustering w/ Gaussians

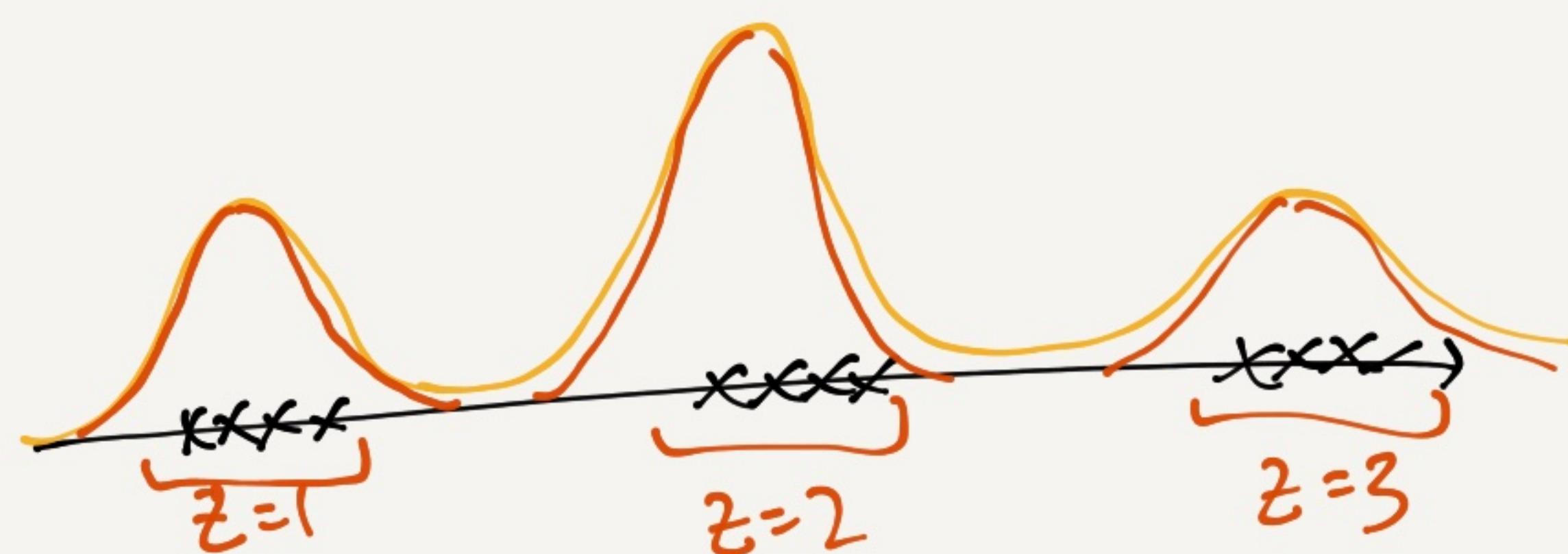
Given data  $D = \{x_1, \dots, x_N\}$

- assume there are  $K$  clusters
- assume each cluster is a Gaussian distribution.

Estimate GMM from  $D$ .

$\Rightarrow$  we obtain the following:

- 1) Cluster center and extent  $(\mu_j, \sigma_j^2)$
- 2) frequency of each cluster  $(\pi_j)$
- 3) the cluster assignment for each  $x_i$  ( $z_i$ )



### Antoni's Hack

- treat  $z_i$  as a parameter to be estimated.
- let  $z_i \in \{1, \dots, K\}$  be the cluster assignment for  $x_i$
- Goal: maximize the joint log-likelihood of  $(x, z)$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_i \log p(x_i, z_i) = \underset{\theta}{\operatorname{argmax}} \sum_i \log p(x_i | z_i) p(z_i)$$

### Indicator variable trick

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

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

$$p(z_i) = \prod_{j=1}^K \pi_j^{z_{ij}} \quad \text{selects the } \pi_{z_i}$$

### Substitute:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log \left[ \prod_{j=1}^K N(x_i | \mu_j, \sigma_j^2)^{z_{ij}} - \prod_{j=1}^K \pi_j^{z_{ij}} \right]$$

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Variables depend on each other, try an alternating maximization scheme:

1) given  $\{\pi_j, \mu_j, \sigma_j^2\}$  fixed, then solve for  $z_{ij}$

- each  $z_{ij}$  is independent of others

$$z_{ij} = \underset{z_{ij}}{\operatorname{argmax}} \sum_i z_{ij} \log \pi_j N(x_i | \mu_j, \sigma_j^2)$$

o, pick j w/ largest  $\log \pi_j N(x_i | \mu_j, \sigma_j^2)$

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

2) given  $z_{ij}$  fixed, solve for  $\{\pi_j, \mu_j, \sigma_j^2\}$

mean  $\mu_j$ :

$$l(\mu_j) = \sum_i -\frac{1}{2} \frac{(x_i - \mu_j)^2}{\sigma_j^2} \cdot z_{ij}$$

$\log N(x_i | \mu_j, \sigma_j^2)$

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

$$= \sum_i z_{ij} x_i - \sum_i z_{ij} \mu_j = 0$$

$$\Rightarrow \boxed{\hat{\mu}_j = \frac{1}{\sum_i z_{ij}} \sum_i z_{ij} x_i}$$

# of points assigned to cluster j.

Sum of points  $x_i$  assigned to j.

mean of points assigned to cluster j.

Similarly

$$\hat{\sigma}_j^2 = \frac{\sum_i z_{ij} (x_i - \hat{\mu}_j)^2}{\sum_i z_{ij}}$$

variance of points in j

$$\hat{\pi}_j = \frac{1}{N} \sum_i z_{ij}$$

] fraction of points assigned to j.

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

Note:

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- this 2-step procedure always maximizes the original objective  $\Rightarrow$  converges to a local maximum

$\bullet \sigma_j^2 = \text{fixed value}, \pi_j = \frac{1}{K}$   
 $\Rightarrow K$ -means clustering algorithm.

- need some initialization of  $(\pi, \mu, \sigma^2)$  or  $(z_e)$  to start the algorithm.  
- different init. may give different answers.

$\bullet$  not maximizing the actual  $p(x)$  (data likelihood), so it's not MLE.

[we are optimizing  $z$ , rather than marginalizing it out as in  $p(x)$ .]

## Expectation Maximization (EM) algorithm

(Dempster, Laird, Rubin 1977 ~ 62K citations on google scholar)  
in 2020.

- MLE when there are hidden variables

$$\begin{cases} X = \text{observation r.v.} \\ Z = \text{hidden r.v.} \\ p(x) = \sum_z p(x|z)p(z) \end{cases}$$

Goal: maximize data LL:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(x) = \underset{\theta}{\operatorname{argmax}} \log \sum_z p(x|z)p(z)$$

Key observation:

- if we had the "complete" data  $\{x, z\}$  observed,

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- then the problem is easy.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(x, z) = \underset{\theta}{\operatorname{argmax}} \sum_i \log p(x_i|z_i) + \log p(z_i)$$

$\Rightarrow$  step 2 in Antonis hack.

- so, guess the values in a probabilistic way.

1) select expected values of  $z$  given previous model  
 $\Rightarrow \hat{z}$

2) maximize  $\log p(x, \hat{z})$  to get new model.

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

Formally

0) Select initial model  $\hat{\theta}^{(\text{old})}$

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

"Q" function  
previous parameters  
new parameters  
we will optimize

conditional expectation  
using the previous model  
 $\hat{\theta}^{(\text{old})}$

2) M-step:  $\hat{\theta}^{(\text{new})} = \underset{\theta}{\operatorname{argmax}} Q(\theta; \hat{\theta}^{(\text{old})})$

3)  $\hat{\theta}^{(\text{old})} \leftarrow \hat{\theta}^{(\text{new})}$ , repeat (1) & (2) until convergence.

EM for GMMs

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

1) E-step:

$$Q(\theta; \hat{\theta}^{(old)}) = E_{z|x, \hat{\theta}^{(old)}} [\log p(x, z)]$$

$$= \sum_i \sum_j E_{z|x, \hat{\theta}^{(old)}} [z_{ij}] \log \pi_j N(x_i | \mu_j, \sigma_j^2)$$

$\hat{z}_{ij}$

(same form as step (2) in my hack, but with  
 $z_{ij}$  replaced with  $\hat{z}_{ij}$ )

$$\hat{z}_{ij} = E_{z|x, \hat{\theta}^{(old)}} [z_{ij}]$$

$$= p(z_{ij}=1 | x, \hat{\theta}^{(old)})$$

Expectation of indicator variable (P.I.J)

$$= \dots \text{(tutorial)}$$

$$= \frac{\hat{\pi}_j N(x_i | \hat{\mu}_j, \hat{\sigma}_j^2)}{\sum_k \hat{\pi}_k N(x_i | \hat{\mu}_k, \hat{\sigma}_k^2)}$$

"soft assignment" to cluster  $j$  using the  $\hat{\theta}^{(old)}$

$$= p(z_{ij}=1 | x_i, \hat{\theta}^{(old)}) \leftarrow \begin{matrix} \text{posterior probability} \\ \text{of } z_{ij} | x_i \end{matrix}$$

Note:  $0 \leq \hat{z}_{ij} \leq 1$   
 soft assignment

$z_{ij} \in \{0, 1\}$   
 hard assignment

M-step: same as before  $z_{ij} \rightarrow \hat{z}_{ij}$

Summary EM-GMM

$$\text{E-step: } \hat{z}_{ij} = p(z_{ij}|x_i, \hat{\theta}^{(old)}) \quad \text{Soft assignment using } \hat{\theta}^{(old)}$$

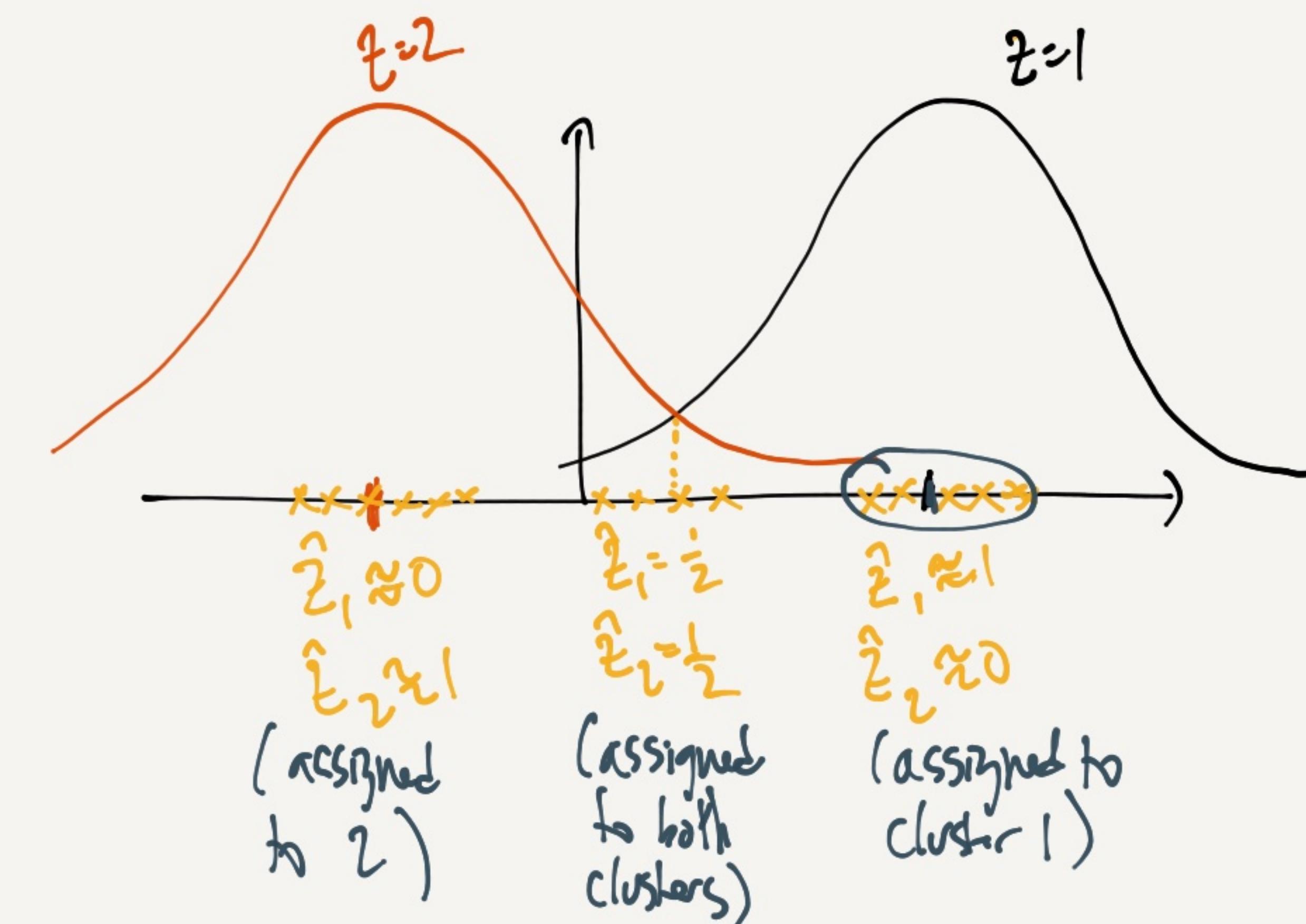
M-Step:

$$\hat{\mu}_j = \frac{1}{\hat{N}_j} \sum_i \hat{z}_{ij} x_i \quad \text{Sample mean w/ weighted points by soft assignments}$$

$$\hat{N}_j = \sum_i \hat{z}_{ij} \quad \text{Total # of points in cluster } j \text{ (soft weighted)}$$

$$\hat{\sigma}_j^2 = \frac{1}{\hat{N}_j} \sum_i \hat{z}_{ij} (x_i - \hat{\mu}_j)^2 \quad "$$

$$\hat{\pi}_j = \frac{1}{N} \hat{N}_j \quad "$$



Assignment Project Exam Help

<https://powcoder.com>

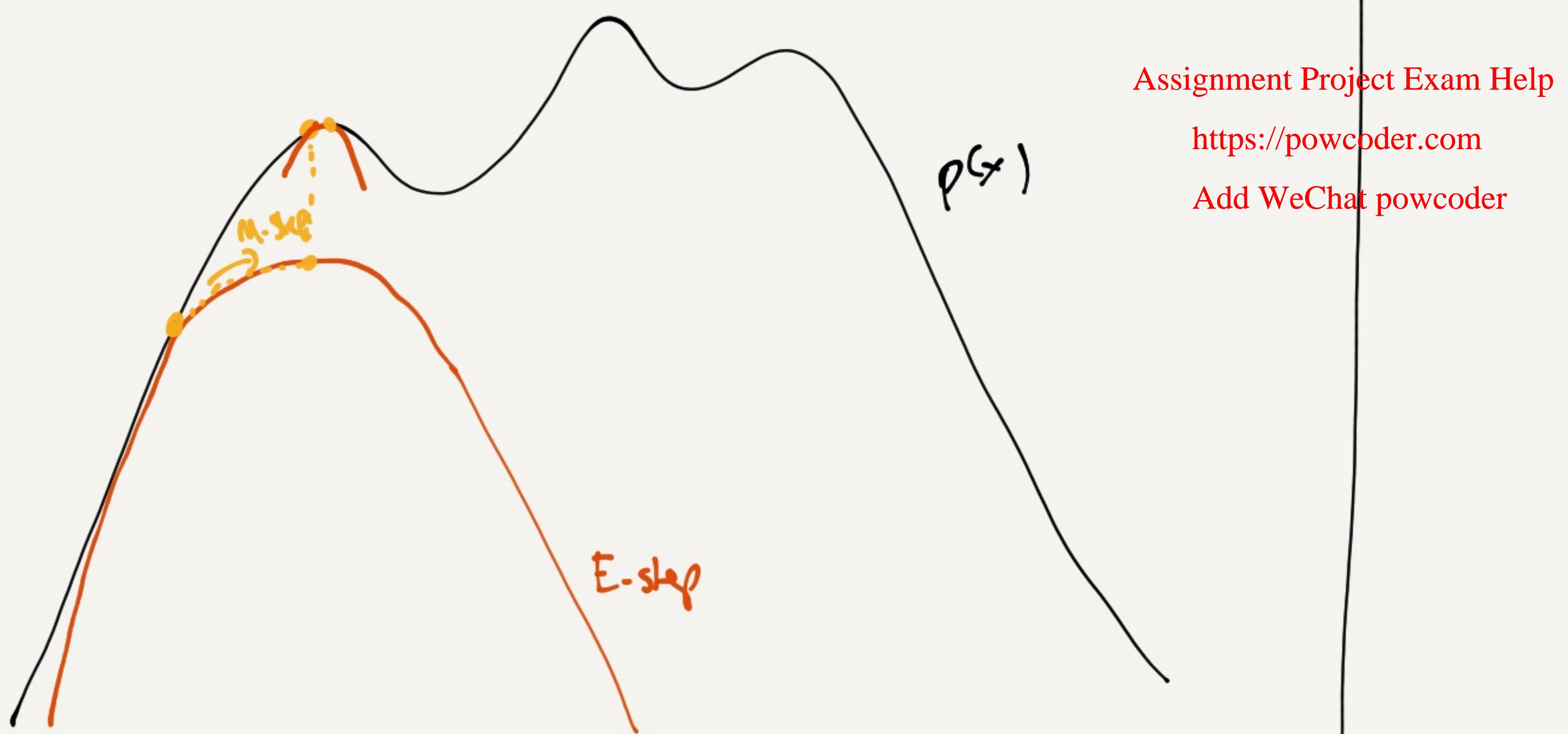
Add WeChat powcoder

## Notes:

1) General algorithm - EM is a general framework for MLE on models w/ hidden variables (not just mixture models)

2) Convergence - after each EM iteration, the data LL increases  $\rightarrow$  converges to local maximum. (could be arbitrarily slow)

3) Initialization - different init give different estimates of  $\theta$ . pick the  $\theta$  w/ maximum  $p(x|\theta)$ .



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Tutorial 4

### 4.12 - Lagrange Multipliers

$$\sum_{j=1}^K \pi_j = 1 \text{ constraint}, \pi_j \geq 0$$

need to solve:  $\hat{\pi}_j = \underset{j}{\operatorname{argmax}} \sum_{j=1}^K z_j \log \pi_j$

$\sum_j \pi_j = 1$   
equality constraint

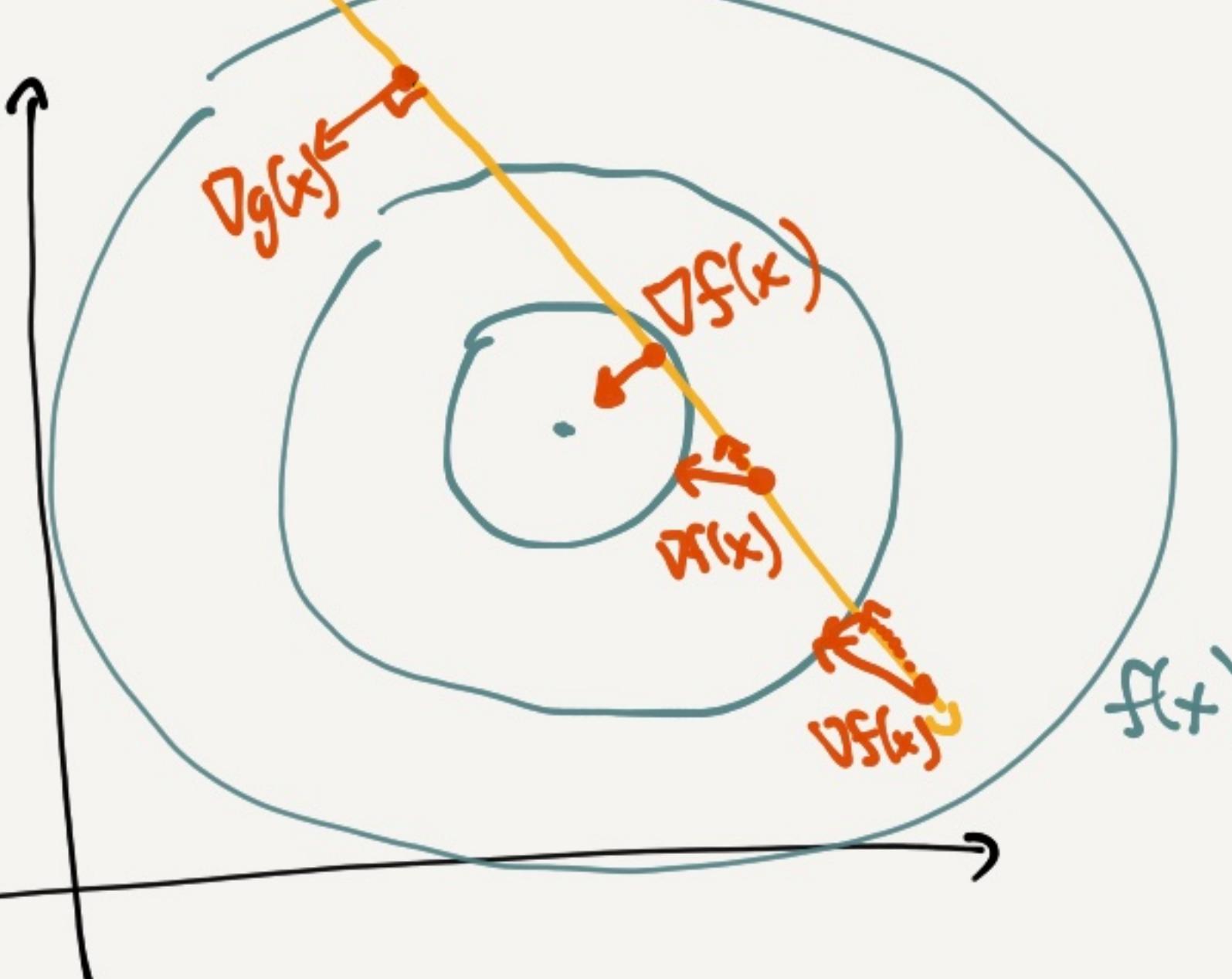
Look at the general problem

$$x^* = \underset{x}{\operatorname{argmax}} f(x) \leftarrow \text{objective function}$$

$$\text{s.t. } g(x) = 0$$

constraint function  
constraint surface

- $\nabla g(x)$  is  $\perp$  to  $g(x) = 0$
- At the optimum,  $\nabla f(x)$  is  $\perp$  to  $g(x) = 0$ .  
(otherwise we could move along  $g(x) = 0$  to increase  $f(x)$ .)



At the optimum, we have

$$\nabla f(x) + \lambda \nabla g(x) = 0, \lambda \neq 0$$

Define Lagrangian function:

$$L(x, \lambda) = f(x) + \lambda g(x)$$

in EM  $\Rightarrow \exists j = N_j$

Find a stationary point  $(x, \lambda)$  of the Lagrangian.

$$\frac{\partial}{\partial x} L(x, \lambda) = 0, \frac{\partial}{\partial \lambda} L(x, \lambda) = 0.$$

a)  $g(\pi) = \sum_{j=1}^K \pi_j - 1 = 0$

$$S(\pi) = \sum_j z_j \log \pi_j$$

Lagrangian:

$$L(\pi, \lambda) = \sum_j z_j \log \pi_j + \lambda \left( \sum_j \pi_j - 1 \right)$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$\frac{\partial L}{\partial \pi_j} = \left( \frac{z_j}{\pi_j} + \lambda \right) \times \pi_j = 0$$

$$\Rightarrow z_j + \lambda \pi_j = 0 \quad \text{run over } j$$

$$\sum_j z_j + \lambda \sum_j \pi_j = 0$$

$$\Rightarrow \lambda = -\frac{\sum z_j}{\sum \pi_j}$$

$$z_j + \left( -\frac{\sum z_j}{\sum \pi_j} \right) \pi_j = 0$$

$$\Rightarrow \boxed{\pi_j = \frac{z_j}{\sum_j z_j}}$$

## 4.6 - mixture of exponentials

$$p(x) = \sum_i \pi_i \lambda_i e^{-x_j x}$$

comp. prob.  $\downarrow$  exponential density  $p(x|j)$

Given data  $D = \{x_i\}_{i=1}^N$   
Define hidden variable  $z_{ij}$  for each  $x_i$ .

Joint LL:

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

E-step:  $Q(\hat{\theta} : \hat{\theta}^{(old)}) = E_{z|x, \hat{\theta}^{(old)}} [\log p(x, z | \theta)]$

$$= E_{z|x, \hat{\theta}^{(old)}} \left[ \sum_i \sum_j z_{ij} \log \pi_j + z_{ij} \log p(x|j) \right]$$

$$= \sum_i \sum_j E_{z|x, \hat{\theta}^{(old)}} [z_{ij}] \log \pi_j + E_{z|x, \hat{\theta}^{(old)}} [z_{ij}] \log p(x|j)$$

$$\hat{z}_{ij} = E_{z|x, \hat{\theta}^{(old)}} [z_{ij}]$$

$\hat{z}_{ij} = E_{z|x, \hat{\theta}^{(old)}} [z_{ij}]$  Indicator variable

$$= p(z_{ij}=1 | x, \hat{\theta}^{(old)}) \cdot 1 + p(z_{ij}=0 | x, \hat{\theta}^{(old)}) \cdot 0$$

$$= p(z_{ij}=1 | x, \hat{\theta}^{(old)})$$
 (expectation of indicator variable is the prob of the indicator = 1)

$$= p(z_{ij}=j | x)$$

$$= \frac{p(x | z_{ij}=j) p(z_{ij}=j)}{p(x)}$$

$$= \frac{p(x_{-i}) p(x_i | z_{ij}=j) p(z_{ij}=j)}{p(x_{-i}) p(x_i)}$$

$$= \frac{p(x_i | z_{ij}=j) p(z_{ij}=j)}{\sum_k p(x_i | z_{ij}=k) p(z_{ij}=k)}$$

$\Rightarrow X = \{x_i, \underbrace{x_1, x_2, \dots}_{X_{-i}}\}$

$x_i \perp x_{-i}$

$\Rightarrow p(x) = p(x_i) \cdot p(x_{-i})$

$$\hat{z}_{ij} = \frac{\pi_j p(x_i | j)}{\sum_k \pi_k p(x_i | k)}$$
 (evaluated w/  $\hat{\theta}^{(old)}$ )

doesn't affect  $x_{-i}$

Note:  $p(x_i | z_{ij}=j) = p(\underbrace{x_{-i}}_{\downarrow} | z_{ij}=j) p(x_i | z_{ij}=j)$

$$= p(x_{-i}) p(x_i | z_{ij}=j)$$

M-step

$$(\hat{\pi}_j, \hat{\lambda}) = \underset{\lambda, \pi}{\operatorname{argmax}} \quad Q(\theta; \hat{\theta}^{(t+1)})$$

$$= \underset{\lambda, \pi}{\operatorname{argmax}} \sum_i \sum_j \hat{z}_{ij} \log \pi_j + \sum_{ij} \log p(x_i | j)$$

Look at  $\pi$  first

$$\hat{\pi}_j = \underset{\pi_j}{\operatorname{argmax}} \sum_i \sum_j \hat{z}_{ij} \log \pi_j$$

$$\sum_j \hat{\pi}_j = 1$$

$\sum_j (\sum_i \hat{z}_{ij}) \log \pi_j$

$\hat{N}_j$

$$= \underset{\text{s.t. } \sum_j \hat{\pi}_j = 1}{\operatorname{argmax}} \sum_j \hat{N}_j \log \pi_j$$

Use 4.12 a)

$$\Rightarrow \hat{\pi}_j = \frac{\hat{N}_j}{\sum_k \hat{N}_k} = \frac{\hat{N}_j}{N}$$

Look at  $\lambda_j$

$$\lambda_j = \underset{\lambda_j}{\operatorname{argmax}} \sum_i \hat{z}_{ij} \log p(x_i | j)$$

$$= \underset{\lambda_j}{\operatorname{argmax}} \sum_i \hat{z}_{ij} (\log \lambda_j - \lambda_j x_i)$$

$$\Rightarrow \frac{\partial}{\partial \lambda_j} = \sum_i \hat{z}_{ij} \left( \frac{1}{\lambda_j} - x_i \right) = 0$$

$$= \frac{1}{\lambda_j} \hat{N}_j - \sum_i \hat{z}_{ij} x_i = 0$$

$$\Rightarrow \frac{1}{\lambda_j} = \frac{1}{\hat{N}_j} \sum_i \hat{z}_{ij} x_i$$

$$\hat{N}_j = \sum_i \hat{z}_{ij}$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$z_i$  = hidden assignment variable  $\in \{1, \dots, k\}$

$z_{ij}$  = indicator that  $z_i=j$  =  $\begin{cases} 1, & z_i=j \\ 0, & \text{otherwise} \end{cases}$

"one-hot-encoding"

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow z_i=1$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \Rightarrow z_i=2$$

$\hat{z}_{ij} = p(z_i=j | x_i)$ , computed in E-step  
, soft assignment of  $x_i$  to  $j$