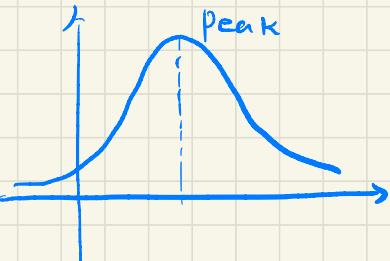
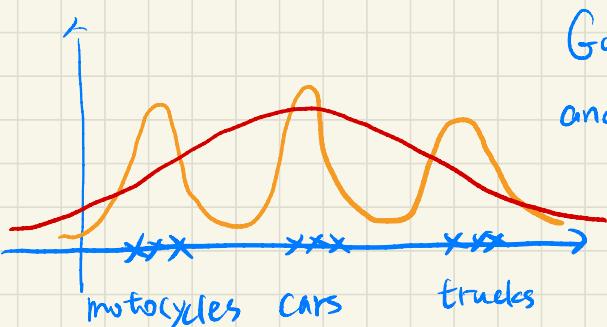


## Lecture 4 Mixture Models and Clustering



so far, we only have looked at probability distribution with respect to one mode (peak)

What if it is more complicated?



Gaussian doesn't fit the data well and doesn't tell the whole story.

### Gaussian Mixture Model

two r.v.

1)  $Z = \text{hidden state}$  (vehicle type) with  $K$  states

e.g.  $Z \in \{\text{scoutter, car, truck}\}$   
1 , 2 , 3

$$P(Z=j) = \pi_j, \quad \sum_j \pi_j = 1$$

(prior probability of a type of vehicle occurring)

2)  $X = \text{observation}$

observation model conditioned on  $Z=j$  (weight)

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

each vehicle type has its own distribution of weight

## Generative process

1) sample  $Z$  (vehicle type)

→ sample  $X|Z$  (weight given type)

Note: We never see  $Z$ ! only see  $X$ !

Distribution of  $X$ :

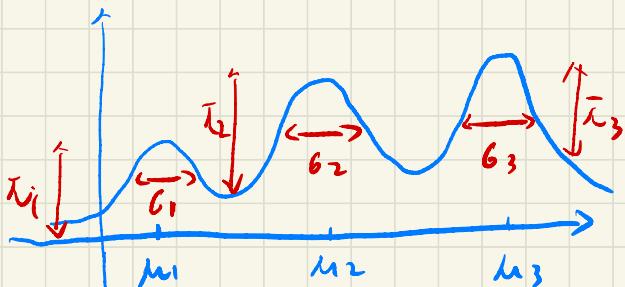
$$p(x) = \sum_j p(x, z=j)$$

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

$$p(x) = \sum_j \pi_j p(x|z=j) \leftarrow \text{"weight sum of Gaussian dist"}$$

$\pi_j$ : component weight (prior)

$p(x|z=j)$ : mixture components



# Clustering

Given data  $D = \{x_1, \dots, x_n\}$ ,  $x_i \in \mathbb{R}$ , estimate a GMM with  $K$  components (# of clusters, Gaussians)

1) Gaussian components  $\mu_j$ ,  $\sigma_j^2$

$\mu_j$ : location of cluster

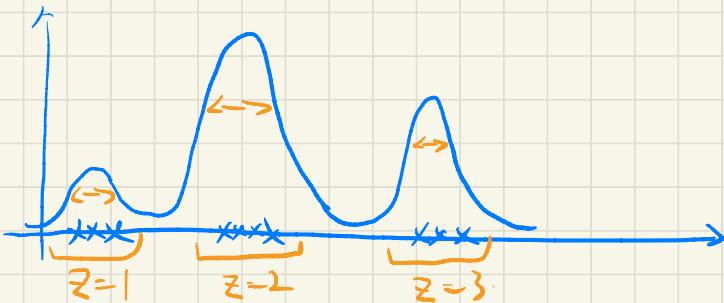
$\sigma_j^2$ : spread of cluster.

2) Component weight  $\pi_j$

$\pi_j$ : probability / size of cluster.

3) Cluster assignments  $z_i$  for each  $x_i$

$z_i$ : cluster membership



Antoni's hack

Data =  $\{x_1, \dots, x_n\}$

Assignment variable  $z_i \in \{1, \dots, k\}$  = cluster assignment for  $x_i$

Objective: treat  $z_i$ 's as a parameter and optimize them  
maximize the joint likelihood  $p(x, z)$

$$(\hat{\theta}, \hat{z}) = \arg \max_{(\theta, z)} \sum_i \log p(x_i, z_i)$$

$$= \arg \max_{(\theta, z)} \sum_i \log p(x_i | z_i) p(z_i)$$

Indicator variable trick:

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

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

$$\Rightarrow \hat{\theta}, \hat{z} = \arg \max \sum_i \left( \log \prod_{j=1}^k \pi_j^{z_{ij}} \cdot \prod_{j=1}^k N(x_i | \mu_j, \sigma_j^2)^{z_{ij}} \right)$$
$$= \arg \max \sum_i \sum_{j=1}^k z_{ij} (\log \pi_j + z_{ij} \log N(x_i | \mu_j, \sigma_j^2))$$

Variables depend on each other, so try an alternative maximization scheme.

1) Given  $\Theta = \{\bar{x}_{ij}, \mu_{ij}, \sigma_{ij}^2\}$ , find the  $Z_{ij}$ 's

- each  $Z_{ij}$  is independent of other  $Z_{ij}$  in the objective

$$\arg \max_{\{Z_{ij}\}} \sum_j Z_{ij} \log \bar{x}_{ij} N(x | \mu_{ij}, \sigma_{ij}^2)$$

↑  
only one term can be selected ( $Z_{ij}=1$ )

$\Rightarrow$  select  $j$  with largest  $\bar{x}_{ij} N(x | \mu_{ij}, \sigma_{ij}^2)$

$$Z_{ij} = \arg \max_j \bar{x}_{ij} N(x_i | \mu_{ij}, \sigma_{ij}^2)$$

2) Given  $Z_{ij}$ , find  $\{\bar{x}_{ij}, \mu_{ij}, \sigma_{ij}^2\}$

$$(\bar{x}_{ij}, \mu_{ij}, \sigma_{ij}^2) = \arg \max_{(\bar{x}_{ij}, \mu_{ij}, \sigma_{ij}^2)} \sum_i Z_{ij} \log \bar{x}_{ij} + Z_{ij} \log N(x_i | \mu_{ij}, \sigma_{ij}^2)$$

mean  $\bar{x}_{ij} = \arg \max_{\mu_{ij}} \sum_i Z_{ij} \left[ -\frac{1}{2\sigma_{ij}^2} (x_i - \mu_{ij})^2 \right]$

$$\frac{\partial}{\partial \mu_{ij}} = \sum_i Z_{ij} \left[ \frac{1}{\sigma_{ij}^2} (x_i - \mu_{ij}) \right] = 0$$

$$= \sum_i Z_{ij} (x_i - \mu_{ij}) = 0$$

$$\Rightarrow \mu_{ij} = \frac{1}{\sum_i Z_{ij}} \left( \sum_i Z_{ij} x_i \right) \leftarrow \begin{matrix} \text{mean of Points} \\ \text{assigned to } j \end{matrix}$$

$\sum_i z_{ij} = \# \text{ of points assigned to cluster } j$

$\sum_i z_{ij} x_i = \text{sum of points assigned to } j$

**Variance**  $\hat{b}_j^2 = \arg \max_{b_j^2} \sum_i z_{ij} \left[ \frac{1}{2} \sqrt{2\pi b_j^2} - \frac{1}{2b_j^2} (x_i - \mu_j)^2 \right]$

$$\begin{aligned} \frac{\partial}{\partial b_j^2} &= \sum_i z_{ij} \left[ \sqrt{2\pi b_j^2} - \frac{1}{2} (2\pi b_j^2)^{-\frac{1}{2}} \cdot 2\pi + \frac{1}{2(b_j^2)^2} (x_i - \mu_j)^2 \right] \\ &= \sum_i z_{ij} \left[ \sqrt{2\pi b_j^2} - \frac{1}{2} (2\pi)^{-\frac{1}{2}} (b_j^2)^{\frac{3}{2}} + \frac{1}{2(b_j^2)^2} (x_i - \mu_j)^2 \right] \\ &= \sum_i z_{ij} \left[ -\frac{1}{2}(b_j^2)^{-1} + \frac{1}{2(b_j^2)^2} (x_i - \mu_j)^2 \right] = 0 \end{aligned}$$

$$\Rightarrow \sum_i z_{ij} \left[ -\frac{1}{2} b_j^2 + \frac{1}{2} (x_i - \mu_j)^2 \right] = 0$$

$$\Rightarrow \hat{b}_j^2 = \frac{\sum_i z_{ij} (x_i - \mu_j)^2}{\sum_i z_{ij}}$$

**Print:**  $\hat{\pi}_j = \arg \max_{\pi_j} \sum_i z_{ij} \log \pi_j + \lambda \left( \sum_{j=1}^K \pi_j - 1 \right)$

$$\frac{\partial}{\partial \pi_j} = \frac{\sum_i z_{ij}}{\pi_j} + \lambda = 0$$

$$\sum_i z_{ij} = -\lambda \pi_j \Rightarrow \hat{\pi}_j = \frac{\sum_i z_{ij}}{-\lambda} = \frac{\sum_i z_{ij}}{N}$$

$$\sum_{j=1}^K \sum_i z_{ij} = -\lambda = N$$

3) Repeat (1) and (2) until converge

- Notes:
- this 2 step procedure always maximizes the objective  $\Rightarrow$  converges to a local maximum
  - need an initial value  $\{z_{ij}\}$  or  $\{\pi_j, \mu_j, b_j^2\}$
  - If we set  $\pi_j = \frac{1}{k}$  and  $b_j^2 = \text{constant}$   
 $\Rightarrow$  k-means algorithm (Lloyd's algorithm)

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

$\mu_j$  = mean of points assigned to  $j$

$$= \frac{1}{\sum_i z_{ij}} \sum_i z_{ij} x_i$$

- problem: not maximizing the actual logP(D)!  
maximizing some surrogate p(X, Z)

Expectation-Maximization (EM) algorithm

Maximum likelihood estimation for models with hidden variables

X = observation r.v.

Z = hidden r.v. r.v.

$$P(X, Z) = P(X|Z) P(Z), P(X) = \sum_Z P(X|Z) P(Z)$$

Goal: MLE

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

key observation

- if we knew  $(X, Z)$ , then problem is easy

$\Rightarrow$  Step 2 of Antoni's hack.

- guess the value of  $Z$  probabilistically :

1) select Expected value of  $Z$  given the model  $\Rightarrow \hat{z}$

2) maximize  $P(x, \hat{z})$  to get the new model

3) repeat 1) and 2)

Formally : EM algorithm

conditional expectation using the

current  $\hat{\theta}^{\text{old}}$

1) select initial model  $\hat{\theta}^{(\text{old})}$

2) E-step :  $Q(\theta; \theta^{(\text{old})}) = \mathbb{E}_{z|x, \hat{\theta}^{(\text{old})}} [\underbrace{\log P(x, z|\theta)}_{\substack{\text{new parameter} \\ \text{old parameter} \\ (\text{fixed})}}]$

Joint LL using  $\theta$

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

4)  $\hat{\theta}^{(\text{old})} \leftarrow \hat{\theta}^{\text{new}}$ , repeat 1 and 2 until converge

$$\text{by } P(D|\theta) = \sum_{i=1}^n \log P(x_i|\theta) = \sum_{i=1}^n \log \sum_{j=1}^K P(x_i, z_j|\theta)$$

assume  $z_j \sim Q(z_j)$

$$= \sum_{i=1}^n \log \sum_{j=1}^K Q(z_j) \frac{P(x_i, z_j|\theta)}{Q(z_j)}$$

$$\geq \sum_{i=1}^n \sum_{j=1}^K Q(z_j) \log \frac{P(x_i, z_j|\theta)}{Q(z_j)} \quad (\text{Jensen})$$

$$\text{IFF: } \frac{P(x_i, z_j|\theta)}{Q(z_j)} = c = \sum_{j=1}^K P(x_i, z_j|\theta)$$

$$\Rightarrow Q(z_j) = \frac{P(x_i, z_j|\theta)}{\sum_{j=1}^K P(x_i, z_j|\theta)} = P(z_j | x_i, \theta)$$

o) Select initial model  $\hat{\theta}^{\text{old}}$ .

$$1) Q(z_j) = \frac{P(x_i, z_j | \hat{\theta}^{\text{old}})}{\sum_{j=1}^K P(x_i, z_j | \hat{\theta}^{\text{old}})}$$

$$M(\theta) = \sum_{i=1}^n \sum_{j=1}^K Q(z_j) \log \frac{P(x_i, z_j | \theta^{\text{old}})}{Q(z_j)} = E(\theta^{\text{old}})$$

$$2) \hat{\theta}^{\text{new}} = \arg \max M(\theta^{\text{old}}) \geq E(\theta^{\text{old}})$$

3) Repeat 1) and 2) until converges.

EM for GMMs

$$\text{Joint LL: } \log p(x, z) = \sum_i \sum_j z_{ij} \log \pi_j N(x_i | \mu_j, \Sigma_j)$$

1) E-step

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

$$= \sum_i \sum_j E_{z|x, \hat{\theta}^{\text{old}}} z_{ij} \underbrace{\log \pi_j N(x_i | \mu_j, \Sigma_j)}_{\text{const.}}$$

$E_{z|x, \hat{\theta}^{\text{old}}} z_{ij}$  } Expectation of an Indicator. PSL-5

$$= P(z_{ij} | x, \hat{\theta}^{\text{old}}) = P(z_i=j | x, \hat{\theta}^{\text{old}})$$

$$= \frac{P(x|z_i=j) P(z_i=j)}{P(x)} \quad \text{Bayes' Rule.}$$

$$= \frac{P(x_{-i}) P(x_i | z_i=j) P(z_i=j)}{P(x_{-i}) P(x_i)} \quad \begin{array}{l} \text{indepent of } x_i \\ \text{w.r.t. oth } x \end{array}$$

$$= \frac{P(x_i | z_i=j) P(z_i=j)}{P(x_i)}$$

$$= \frac{\pi_j N(x_i | \mu_j, \Sigma_j)}{\sum_j \pi_j N(x_i | \mu_j, \Sigma_j)} \quad \begin{array}{l} \text{"soft assignment to} \\ \text{cluster } j \text{ using } \hat{\theta}^{\text{old}} \end{array}$$

$$= P(Z_i=j | x_i, \hat{\theta}^{old}) \quad \text{posterior prob of } Z_i=j \\ \text{(using } \hat{\theta}^{old} \text{)}$$

2) M-step: Same as befr, replace  $Z_{ij}$  with  $\hat{Z}_{ij}$

Summary Em-Gauss

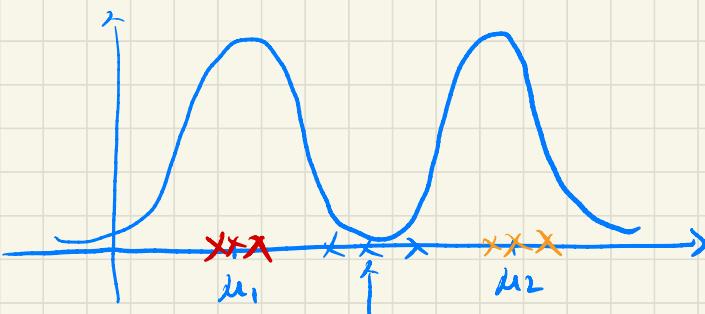
$$\text{E-step: } \hat{Z}_{ij} = P(Z_{ij} | x_i, \hat{\theta}^{old}) = \frac{\pi_j N(x_i | \mu_j, \sigma_j^2)}{\sum_k \pi_k N(x_i | \mu_k, \sigma_k^2)} \\ \text{(using } \hat{\theta}^{old} \text{)}$$

$$\text{M-step: } \hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{Z}_{ij} x_i \leftarrow \begin{array}{l} \text{sample mean with points} \\ \text{weighted by soft assignment } \hat{Z}_{ij} \end{array}$$

$$N_j = \sum_i \hat{Z}_{ij} \leftarrow \text{weight of points assigned to } j$$

$$\hat{\sigma}_j^2 = \frac{1}{N_j} \sum_i \hat{Z}_{ij} (x_i - \hat{\mu}_j)^2$$

$$\hat{\pi}_j = N_j / N$$



$$\hat{z}_{i1} = 1/2$$

$$\hat{z}_{i2} = 1/2$$

## Notes on EM:

1) converges: after each iteration of EM, the data LL increases  $\rightarrow$  converges to local max and could be slow

2) depends on initialization

different init  $\rightarrow$  different  $\hat{\theta}$

pick  $\hat{\theta}$  with largest LL  $p(x|\hat{\theta})$

3) general framework for MLE on any model

with hidden variables: Linear dynamical system

Hidden Markov model

prob. graphical models.