# Lecture 4 - Mixture Models & Clustering
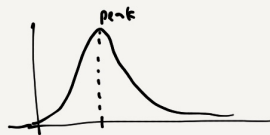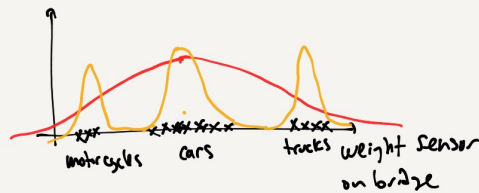


So far, we only have looked at prob. dist. w/ one mode (peak).

## What if it's more complicated?



motorcycles   cars   trucks   weight sensor on bridge

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

## Gaussian Mixture Model

**two r.v.** ← with K states.

(1) $Z$ = hidden state (vehicle type)

e.g. $Z \in \{ \underset{1}{scooter}, \underset{2}{car}, \underset{3}{truck} \}$

$p(Z=j) = \pi_j$ , $\sum_j \pi_j = 1$

$\underbrace{\phantom{xxxx}}$ prior probability of a type of vehicle occurring.

(2) $X$ = observation. observation model conditioned on $Z=j$ (weight)

$$p(x \mid Z=j) = N(x \mid \mu_j, \sigma_j^2)$$

each vehicle type has its own distribution of weight

## Generative Process

1) sample $Z$ (vehicle type)

2) sample $X \mid Z$ (weight given type)

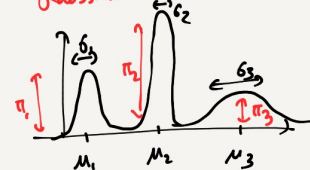Note: we never see $Z$! only see $x$!

Distribution of $x$:

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

$$= \sum_j p(x \mid Z=j) \, p(Z=j)$$

$$\boxed{p(x) = \sum_j \pi_j \, N(x \mid \mu_j, \sigma_j^2)}$$ ← "weighted sum of Gaussian distributions"

component weight (prior)
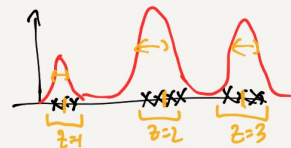
mixture components



## Clustering

$x_i \in \mathbb{R}$

Given data $D = \{x_0, \ldots, x_n\}$, estimate a GMM w/ K components

⟹ 1) Gaussian component $\mu_j, \sigma_j^2$

location of cluster

spread of cluster

# of clusters (Gaussians)

2) component weight $\pi_j$ ← probability / size of cluster

3) cluster assignments $z_i$ for each $x_i$ (cluster membership)

| Antoni's hack |
| --- |

Data   $D = \{x_1, ..., x_n\}$

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

Objective: treat $z_i$'s as a parameter, and optimize them.

maximize the joint likelihood $p(x, z)$:

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

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

Indicator variable trick

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

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

($z_{ij}$ selects the $\pi_j$ & $(\mu_j, \sigma_j^2)$ for the $x_i$)

$$\Rightarrow \hat{\theta}, \hat{z} = \text{argmax} \sum_i \log \prod_j \pi_j^{z_{ij}} \prod_j N(x_i | \mu_j, \sigma_j^2)^{z_{ij}}$$

$$\hat{\theta}, \hat{z} = \underset{\theta, z}{\text{argmax}} \sum_i \sum_j z_{ij} \log \pi_j + z_{ij} \log N(x_i | \mu_j, \sigma_j^2)$$

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

1) Given $\theta = \{\pi_j, \mu_j, \sigma_j^2\}$, find the $z_i$'s.
   - each $z_i$ is independent of other $z_i$s in the objective.

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

⇐ only 1 term can be selected ($z_{ij} = 1$)

$$\Rightarrow \text{select } j \text{ w/ largest } \pi_j N(x_i | \mu_j, \sigma_j^2)$$

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

2) Given $z_i$, find $\{\pi_j, \mu_j, \sigma_j^2\}$

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

Mean $\hat{\mu}_j = \underset{\mu_j}{\text{argmax}} \sum_i z_{ij} \left( -\frac{1}{2\sigma_j^2} (x_i - \mu_j)^2 \right)$

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

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

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

← sum of points assigned to j.

$$\Rightarrow \boxed{\mu_j = \frac{1}{\sum_i z_{ij}} \sum_i z_{ij} x_i}$$

← # of points assigned to cluster j.

← mean of points assigned to j.

Similarly,

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

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

← fraction of points assigned to j

← variance of points assigned to j.

3) Repeat (1) & (2) until convergence

Notes: • this 2-step procedure always maximizes the objective
$\Rightarrow$ converges to a local maximum.

• need an initial value $\{z_i\}$ or $\{\pi_j, \mu_j, \sigma_j^2\}$

• if we set $\pi_j = \frac{1}{K}$ & $\sigma_j^2 =$ constant,
$\Rightarrow$ K-means algorithm (Lloyd's algorithm)

$\begin{cases} z_i = \underset{j}{\arg\min}\ (x_i - \mu_j)^2 \\[2mm] \mu_j = \text{mean of points assigned to } j \\[1mm] \qquad = \frac{1}{\sum_i z_{ij}} \sum_i z_{ij} x_i \end{cases}$

• problem: <u>not</u> maximizing the actual $\log p(D)$ !
maximizing some surrogate $p(x, z)$.

---

Expectation - Maximization (EM) algorithm

(Dempster, Laird, Rubin) 1977 $\Rightarrow$ 66,000 citations on Google Scholar.

Maximum likelihood estimation for models w/ hidden variables
$X =$ observation r.v.
$Z =$ hidden r.v.
$p(X, Z) = p(X|Z)\,p(Z)$ , $\underline{p(X) = \sum_Z p(X|Z)\,p(Z)}$

<u>Goal:</u> <u>MLE</u>
$\theta = \underset{\theta}{\arg\max}\ \log p(X) = \underset{\theta}{\arg\max}\ \log \sum_Z p(X|Z)\,p(Z)$

<u>Key Observation</u>

– if we knew $(X, Z)$, then the problem is easy
$\Rightarrow$ Step 2 of Antoni's hack

– guess the value of $Z$ probabilistically:
i) select Expected value of $Z$ given the model $\Rightarrow \hat{z}$
ii) maximize $p(X, \hat{z})$ to get the new model
iii) repeat 1 & 2.

<u>Formally:</u> EM algorithm

0) Select initial model $\hat{\theta}^{(old)}$
1) <u>E-step:</u> $Q(\theta; \theta^{(old)}) = E_{Z|X,\ \hat{\theta}^{(old)}}\left[\log p(X, Z|\theta)\right]$

<span style="color:red">joint LL using $\theta$</span>

<span style="color:red">new param</span> <span style="color:red">old params</span> <span style="color:red">conditional expectation (fixed) using the current $\hat{\theta}^{(old)}$</span>

2) <u>M-step:</u> $\hat{\theta}^{new} = \underset{\theta}{\arg\max}\ Q(\theta; \hat{\theta}^{(old)})$
3) $\hat{\theta}^{(old)} \leftarrow \hat{\theta}^{new}$ , repeat 1 & 2 until convergence.

# EM for GMMs

Joint LL: $\log p(X,Z) = \sum_i \sum_j z_{ij} \log \pi_j N(x_i | \mu_j, \sigma_j^2)$

## 1) E-step

$$Q(\theta; \hat{\theta}^{old}) = E_{Z|X, \hat{\theta}^{old}}\left[\log p(X,Z)\right]$$

$$= \sum_i \sum_j E_{Z|X}[z_{ij}] \log \pi_j N(x_i | \mu_j, \sigma_j^2) \quad \Leftarrow \text{Same form as in Antoni's hack}$$

$\underbrace{\qquad}_{\hat{z}_{ij}}$

$\hat{z}_{ij} = E_{Z|X, \hat{\theta}^{old}}[z_{ij}]$  → Expectation of an indicator PI-5

$= p(z_{ij} | X, \hat{\theta}^{old}) = p(z_i = j | X, \hat{\theta}^{old})$  → Bayes Rule

$= \dfrac{p(X | z_i = j) p(z_i = j)}{p(X)}$

→ independence of $x_i$ w.r.t. other $X$

$p(X) = p(x_i) p(X_{\setminus i})$

not i (the other $X$'s)

$= \dfrac{p(X_{\setminus i}) p(x_i | z_i = j) p(z_i = j)}{p(X_{\setminus i}) p(x_i)}$

$\left[\hat{z}_{ij} = \dfrac{\hat{\pi}_j N(x_i | \hat{\mu}_j, \hat{\sigma}_j^2)}{\sum_k \hat{\pi}_k N(x_i | \hat{\mu}_k, \hat{\sigma}_k^2)}\right\}$ "soft assignment" to cluster $j$ using $\hat{\theta}^{old}$.

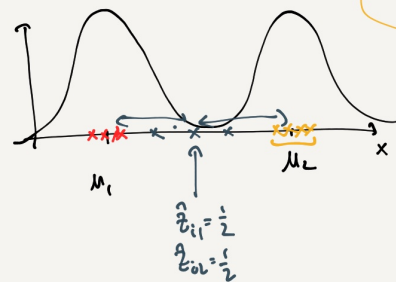$= p(z_i = j | x_i, \hat{\theta}^{old})$  ← posterior prob of $z_i | x_i$ (using $\hat{\theta}^{old}$)

## 2) M-step: same as before, replace $z_{ij}$ w/ $\hat{z}_{ij}$

---

# Summary EM-GMM

E-step: $\hat{z}_{ij} = p(z_i = j | x_i) = \dfrac{\pi_j N(x_i | \mu_j, \sigma_j^2)}{\sum_k \pi_k N(x_i | \mu_k, \sigma_k^2)}$ ] using $\hat{\theta}^{old}$

M-step:
$\hat{\mu}_j = \dfrac{1}{N_j} \sum_i \hat{z}_{ij} x_i$ ← sample mean w/ points weighted by soft assignment $\hat{z}_{ij}$.

$N_j = \sum_i \hat{z}_{ij}$ ← weight of points assigned to $j$

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

$\hat{\pi}_j = N_j / N$ ←

w/ = with
w/o = without
wrt = with respect to



$\mu_1 \qquad \mu_2 \qquad x$

$\hat{z}_{i1} = \frac{1}{2}$
$\hat{z}_{i2} = \frac{1}{2}$

# Notes on EM:

1) <u>converges</u> – after each iteration of EM, the data LL
   increases → converges to a local max.
   (could be slow)

2) <u>depends on initialization</u>
   different init → different $\hat{\Theta}$
   pick $\hat{\Theta}$ w/ largest LL $p(x|\hat{\Theta})$

3) <u>general framework</u> for MLE on any model
   w/ hidden variables: Linear dynamical system
   Hidden Markov model
   prob. graphical models ...