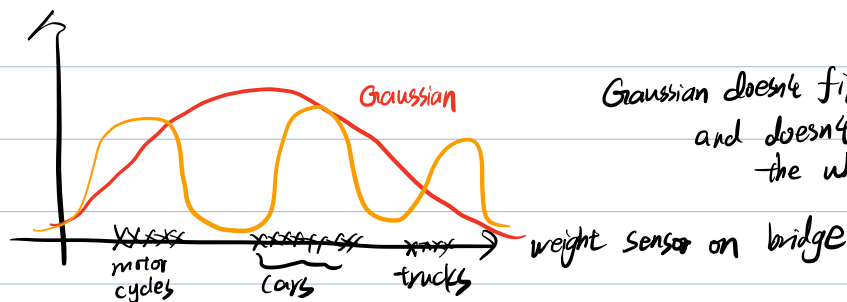


So far, we only have models / has prob. dist.
w/ one mode \Rightarrow has 1 peak

What if it's more complicated?



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

Gaussian Mixture Model (GMM)

two r.v.

1) z = hidden state. (vehicle type in example)

z with k states.

eg. $z \in \{ \text{scooter}_1, \text{car}_2, \text{truck}_3 \}$

$$p(z=j) = \pi_j, \quad \sum \pi_j = 1$$

\uparrow prior probability of type of vehicle occurring

2) x = observation

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

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

$\nearrow \nearrow$
each vehicle type has its own
distribution of weight.

Generative Process

1) sample z (vehicle type)

2) sample $x|z$ (weight given type)

Note: we never see z ! only x (observation)

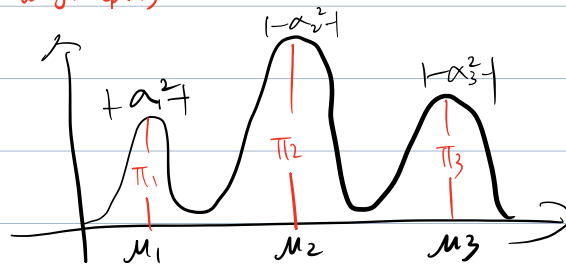
Distribution of X :

$$p(x) = \sum_j p(x, z=j) \text{ (marginalize all } z)$$

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

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

π_j : component weight (prior)
 $\mathcal{N}(x | \mu_j, \sigma_j^2)$: mixture components
 2: weighted sum of Gaussian distributions



Clustering (XG-R)

Given data $D = \{x_1, \dots, x_n\}$, estimate a GMM w/
K components.

⇒ 1D Gaussian components μ_j, σ_j^2
 \uparrow \nwarrow spread of cluster
 location of cluster

2) component weight $\pi_j \propto \text{probability/size of cluster}$

37 cluster assignments z_i for each x_i \hookrightarrow cluster membership.

Antoni's hack

Data $D = \{x_1, x_2, \dots, x_n\}$

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

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

$$\hat{\theta}, \hat{z} = \underset{\theta, z}{\operatorname{argmax}} \sum_i \log p(x_i, z_i) \quad \text{joint pdf of } (x_i, z_i)$$

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

Indicator variable trick

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

$$p(z_i) = \prod_{j=1}^K \pi_j^{z_{ij}} \quad (\text{like bernoulli distribution})$$

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

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

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

(1) Given $\theta = \{\pi_j, \mu_j, \sigma_j^2\}_j$, find z_i

- each z_i is independent of others.

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

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

(in other words, $z_i = \underset{j}{\operatorname{argmax}} \pi_j \mathcal{N}(x_i | \mu_j, \sigma_j^2)$)

since only one
of $z_{ij} = 1$

(2) Given z_i , find $\theta = \{\pi_j, \mu_j, \sigma_j^2\}_j$

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

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

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

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

sum of points
assigned to j

\rightarrow # of points assigned to
cluster j

the whole thing \Rightarrow mean of points assigned to j

Similarly, $\sigma_j^2 = \frac{\sum_i z_{ij} (x_i - \mu_j)^2}{\sum_i z_{ij}}$ \leftarrow variance of points assigned to j

$\pi_j = \frac{\sum_i z_{ij}}{N}$ \leftarrow fraction of points assigned to j

23) • repeat (1) and (2) until it converges

\Rightarrow converge to a local maximum

- need an initial value $\{z_i\}_i$ or $\{\pi_j, \mu_j, \sigma_j^2\}_j$
- if we set $\pi_j = \frac{1}{K}$ & $\sigma_j^2 = \text{constant}$

\Rightarrow K-means

$$\left\{ \begin{array}{l} z_i = \arg\min_j (x_i - \mu_j)^2 \end{array} \right.$$

$$\begin{aligned} \mu_j &= \text{mean of points assigned to } j \\ &= \frac{1}{\sum_i z_{ij}} \cdot \sum_i z_{ij} x_i \end{aligned}$$

- problem: not maximizing the actual $\log p(D)$!

maximizing some surrogate $p(X, Z)$

Z is r.v., but we treat it as a parameter.

Expectation - Maximization (EM) algorithm (create MLE of above case)

(Dempster, Laird, Rubin) 1977 \Rightarrow 66000 citations

Maximum likelihood estimation for models
with hidden variables

X = observation r.v.

Z = hidden r.v.

$$p(X, Z) = p(X|Z) \cdot p(Z) \quad , \quad p(X) = \sum_Z p(X|Z) p(Z)$$

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

\downarrow
tricky $\Rightarrow Z$ inside log.

Key observation.

- if we know (X, Z) , then the 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 1-2

Formally

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

1) E-step: $Q(\theta, \hat{\theta}^{(old)})$

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