

04/27/15

(1)

# UNSUPERVISED LEARNING

"Make sense of unlabelled data"

① K-Means: "Show me different types of input in my dataset"

Given: unlabelled data  $D = \{ \vec{x}_1, \vec{x}_2, \dots, \vec{x}_N \}$   
to some fixed integer  $k$

Goal: Find  $k$  "clusters" or groupings of data

Variables

- Cluster Centers:  $\vec{\mu}_1, \vec{\mu}_2, \dots, \vec{\mu}_k \in \mathbb{R}^d$
- Cluster-Assignment Ids:  $c_1, c_2, \dots, c_N \in \{1, 2, \dots, k\}$
- Cluster-Assignment:  $\vec{a}_1, \vec{a}_2, \dots, \vec{a}_k \in \{0, 1\}^k$   
1-of- $k$  encoding

"assignment" or "association" vector for pt  $i$

$\vec{a}_i = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  ← 4<sup>th</sup> cluster

$a_{ij} = 1 \iff$  datapoint  $i$  belongs to cluster  $j$   
 $a_{ij} \in \{0, 1\} \quad \sum_{j=1}^k a_{ij} = 1$

K-Means Algorithm:

→ Assignment Step:

$$c_i \leftarrow \underset{j}{\operatorname{argmin}} \|\vec{x}_i - \vec{\mu}_j\|_2^2$$

Assign each data-pt to closest cluster center; closest in  $L_2$ -sense

→ Recenter Step

$$\mu_j \leftarrow \frac{1}{\sum_{i=1}^N a_{ij}} \sum_{i=1}^N \vec{x}_i$$

#pts-assigned-to-j

⇔  $\mu_j$  is centroid/mean of  $\vec{x}_i$  assigned to  $j$

② K-means as an Optimization problem/alg.

- So far, we have a "procedural" view of K-means
- But what is K-means doing?
- Does it ever converge?
- What does it all mean???

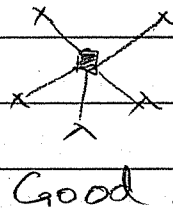
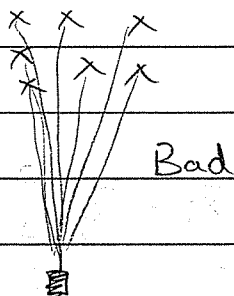
Objective:

$$F(\vec{\mu}_1, \dots, \vec{\mu}_k, C_1, \dots, C_k) = \sum_{i=1}^N \|\vec{x}_i - \vec{\mu}_{c_i}\|_2^2$$

$$F(\vec{\mu}_1, \dots, \vec{\mu}_k, \vec{a}_1, \dots, \vec{a}_k) = \sum_{i=1}^N \sum_{j=1}^k a_{ij} \|\vec{x}_i - \mu_j\|_2^2$$

"Distortion Function": sum of squared distances between pts & cluster-centers-assigned-to-them

Picture:



Under this view / objective, the optimal clustering:

$$\min_{\vec{\mu}_1, \dots, \vec{\mu}_k} \min_{\vec{a}_1, \dots, \vec{a}_N} F(\vec{\mu}, \vec{a})$$

Mixed Continuous / Discrete Optimisation  
Very very hard!

Approximate optimization Alg: Co-ordinate Descent

→ Fix  $\{\vec{\mu}_j\}$  Optimize  $\{a_{ij}\}$  (or  $\{C_i\}$ )

$$\min_{C_1, \dots, C_N} \sum_{i=1}^N \|\vec{x}_i - \vec{\mu}_{C_i}\|^2$$

$$= \sum_{i=1}^N \min_{C_i} \|\vec{x}_i - \vec{\mu}_{C_i}\|^2 \equiv \text{Assignment Step}$$

Set  $C_i = j^*$  that minimizes distance (Closest center center)

→ Fix  $\{C_i\}$  (or  $\{a_{ij}\}$ ), min  $\{\vec{\mu}_j\}$

$$\min_{\vec{\mu}_1, \dots, \vec{\mu}_k} \sum_{i=1}^N \sum_{j=1}^k a_{ij} \|\vec{x}_i - \vec{\mu}_j\|^2$$

↔ swap

$$= \min_{\vec{\mu}_1, \dots, \vec{\mu}_k} \sum_{j=1}^k \sum_{i=1}^N a_{ij} \|\vec{x}_i - \vec{\mu}_j\|^2$$

$$= \sum_{j=1}^k \min_{\vec{\mu}_j} \sum_{i: a_{ij}=1} \|\vec{\mu}_j - \vec{x}_i\|^2$$

Centroid Problem = Re-center Step  
= Set  $\vec{\mu}_j = \text{mean of } \vec{x}_i$  st  $a_{ij} = 1$

### ③ Gaussian Mixture Model (GMM)

"Tell me how likely or rare each data point is"

Given: Data  $D = \{\vec{x}_1, \dots, \vec{x}_N\}$   $\vec{x}_i \in \mathbb{R}^d$

Goal:  $\vec{x} \sim P(\vec{x})$

"Fit"  $P(\vec{x})$  to dataset

Note:  $P(\vec{x})$  is not Gaussian (multimodal data)

But all I know are Gaussians "!"

Maybe it's a "collection of Gaussians?"

Formally, we introduce a "Latent" or "Hidden" variable  $Z = z \in \{1, 2, \dots, k\}$  that indexes  $k$  Gaussians

So GMM:

$$\textcircled{1} Z \sim \text{Cat}(\vec{\pi}) \Leftrightarrow P(Z=c) = \pi_c \quad \begin{array}{l} \pi_c \geq 0 \\ \sum_{c=1}^k \pi_c = 1 \end{array}$$

$$\textcircled{2} \vec{x} | Z=c \sim N(\vec{\mu}_c, \Sigma_c) \quad \left. \begin{array}{l} \text{ie } \vec{x} \text{ given } Z=c \text{ is a} \\ \text{Gaussian w/ mean vector } \vec{\mu}_c \in \mathbb{R}^d \\ \text{covariance matrix } \Sigma_c \in \mathbb{R}^{d \times d} \end{array} \right\}$$

$$P(\vec{x}=\vec{x} | Z=c) = \frac{1}{(2\pi)^{d/2} |\Sigma_c|^{1/2}} e^{-\frac{(\vec{x}-\vec{\mu}_c)^T \Sigma_c^{-1} (\vec{x}-\vec{\mu}_c)}{2}}$$

$$\text{So } P(\vec{x}=\vec{x}) = \sum_{c=1}^k P(\vec{x}=\vec{x}, Z=c)$$

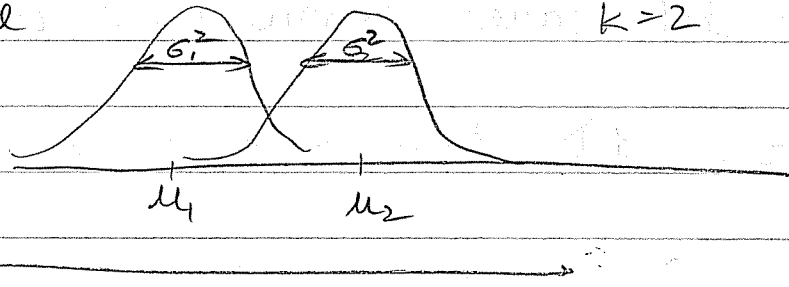
$$= \sum_{c=1}^k P(\vec{x}=\vec{x} | Z=c) P(Z=c)$$

$$\text{mixture} = \sum_{c=1}^k \pi_c \underbrace{N(\vec{\mu}_c, \Sigma_c)}_{\text{Gaussian}}$$

3

$d = 1$   
 $k = 2$

1D example



### ④ Parameter learning in GMMs

Given:  $D$

Estimate Parameters:  $\{\pi_1, \dots, \pi_k, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k\} \equiv \Theta$

How? How else? Max-likelihood! (Max-Marginal-likelihood actually)

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \prod_{i=1}^N P(\vec{X} = \vec{x}_i | \Theta)$$

$$\equiv \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^N \log P(\vec{X} = \vec{x}_i | \Theta)$$

Problem!  $Z$  not available at training time

Must marginalize  $\Leftrightarrow$  sum inside log

$\Leftrightarrow$  all parameters coupled

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^N \log \sum_{c=1}^k P(\vec{X} = \vec{x}_i, Z=c | \Theta)$$

$$= \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^N \log \sum_{c=1}^k P(\vec{X} = \vec{x}_i | Z=c, \Theta) P(Z=c | \Theta)$$

Doesn't "factorize"

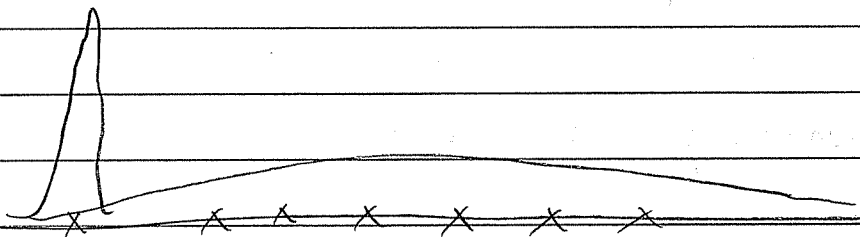
Compare to supervised learning

$$D = \{(\vec{x}_1, z_1), (\vec{x}_2, z_2) \dots (\vec{x}_N, z_N)\}$$

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^N \left[ \underbrace{\log P(\vec{X} = \vec{x}_i | Z = z_i, \Theta)}_{\text{Separate}} + \underbrace{\log P(Z=c | \Theta)}_{\text{Separate}} \right]$$

⑤ Hidden data causes singularities in marginal likelihood

Consider 1D data with  $k=2$



What if I "fit" a very very peaky Gaussian to 1 point (say  $x_1$ ) to another Gaussian to everything else.

So  $\vec{\mu}_1 = \vec{x}_1$

$$P(X=x_1 | Z=1) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_1 - \mu_1)^2}{2\sigma^2}}$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^0$$

So if  $\sigma_1 \rightarrow 0$   $P(x_1 | Z=1, \vec{\theta}) \rightarrow \infty$   
 $\Rightarrow$  likelihood  $\rightarrow \infty$

(4)

## ⑥ Relating Max- (Marginal) Likelihood Estimation M(M)LE in GMMs to K-Means

Assumption ①  $\pi_c = \frac{1}{K}$  All Gaussians/Clusters equally likely a priori

Assumption ② Same Spherical Gaussian  $\forall c$

$$\Sigma_c = \Sigma = \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}_{d \times d} = \sigma^2 I_{d \times d}$$

$$\Leftrightarrow P(\vec{X} = \vec{x} | Z = c) \propto e^{-\frac{1}{2}(\vec{x} - \mu_c)^T \frac{1}{\sigma^2} (\vec{x} - \mu_c)}$$

$$= e^{-\frac{\|\vec{x} - \mu_c\|^2}{2\sigma^2}}$$

Assumption ③ "Hard Assignment": Each  $\vec{x}_i$  "belongs to" exactly 1 cluster / Gaussian

$$P(\vec{X} = \vec{x}, Z = c) \propto \begin{cases} 0 & \text{if } z=1 \\ 0 & \text{if } z=2 \\ P(\vec{X} = \vec{x} | Z = c) \\ \vdots \\ 0 \end{cases}$$

Assumptions ①, ②, ③

$$\Rightarrow \text{Marginal likelihood} = \sum_{i=1}^N \log P(\vec{X} = \vec{x}_i | \Theta) = \sum_{i=1}^N \log \sum_c P(\vec{x}_i, c | \Theta)$$

$$= - \sum_{i=1}^N \frac{\|\vec{x}_i - \mu_{c_i}\|^2}{2\sigma^2} + \text{const}$$

$\Rightarrow$  M(M)LE  $\equiv$  K-Means!

1.  $\int \frac{1}{x^2} dx = \int x^{-2} dx = \frac{x^{-1}}{-1} + C = -\frac{1}{x} + C$

2.  $\int \frac{1}{x^3} dx = \int x^{-3} dx = \frac{x^{-2}}{-2} + C = -\frac{1}{2x^2} + C$

3.  $\int \frac{1}{x^4} dx = \int x^{-4} dx = \frac{x^{-3}}{-3} + C = -\frac{1}{3x^3} + C$

4.  $\int \frac{1}{x^5} dx = \int x^{-5} dx = \frac{x^{-4}}{-4} + C = -\frac{1}{4x^4} + C$

5.  $\int \frac{1}{x^6} dx = \int x^{-6} dx = \frac{x^{-5}}{-5} + C = -\frac{1}{5x^5} + C$

6.  $\int \frac{1}{x^7} dx = \int x^{-7} dx = \frac{x^{-6}}{-6} + C = -\frac{1}{6x^6} + C$

7.  $\int \frac{1}{x^8} dx = \int x^{-8} dx = \frac{x^{-7}}{-7} + C = -\frac{1}{7x^7} + C$

8.  $\int \frac{1}{x^9} dx = \int x^{-9} dx = \frac{x^{-8}}{-8} + C = -\frac{1}{8x^8} + C$

9.  $\int \frac{1}{x^{10}} dx = \int x^{-10} dx = \frac{x^{-9}}{-9} + C = -\frac{1}{9x^9} + C$

10.  $\int \frac{1}{x^{11}} dx = \int x^{-11} dx = \frac{x^{-10}}{-10} + C = -\frac{1}{10x^{10}} + C$

11.  $\int \frac{1}{x^{12}} dx = \int x^{-12} dx = \frac{x^{-11}}{-11} + C = -\frac{1}{11x^{11}} + C$

12.  $\int \frac{1}{x^{13}} dx = \int x^{-13} dx = \frac{x^{-12}}{-12} + C = -\frac{1}{12x^{12}} + C$