

Last time: Boosting

$$\text{Ensemble: } h_m(x) = \sum_{j=1}^m \hat{\alpha}_j \underbrace{h(x; \theta_j)}_{\text{base learner}} \xrightarrow{\text{Vote}}$$

$$\text{Parameters } \Theta = \{(\hat{\alpha}_j, \hat{\theta}_j)\}_{j=1}^m \Rightarrow \hat{\theta}_j = \{s, \theta_0, \underbrace{\theta_k}_{k}\}.$$

This is a mixture model and we want to find more general ones.

Today:

- Gaussian Mixture Model

- Inference / Learning of Parameters in GMM

given  $\left\{ \begin{array}{l} \text{incomplete } (\text{EM algorithm}) \\ \text{complete} \end{array} \right.$

Basic Mixture model :  $\underbrace{\mathcal{D}}_{\text{also denote as } X} = \{x_t\}_{t=1}^n \Rightarrow \text{unsupervised}$

Capture and resolve observable ambiguities in data

→ m-component GMM is :

$$P(x; \Theta) = \sum_{j=1}^m P(j) N(x; \mu_j, \Sigma_j)$$

①  $\{P(j)\}_{j=1}^m$  : mixing parameters  $P(j) \geq 0 \& \sum P(j) = 1$

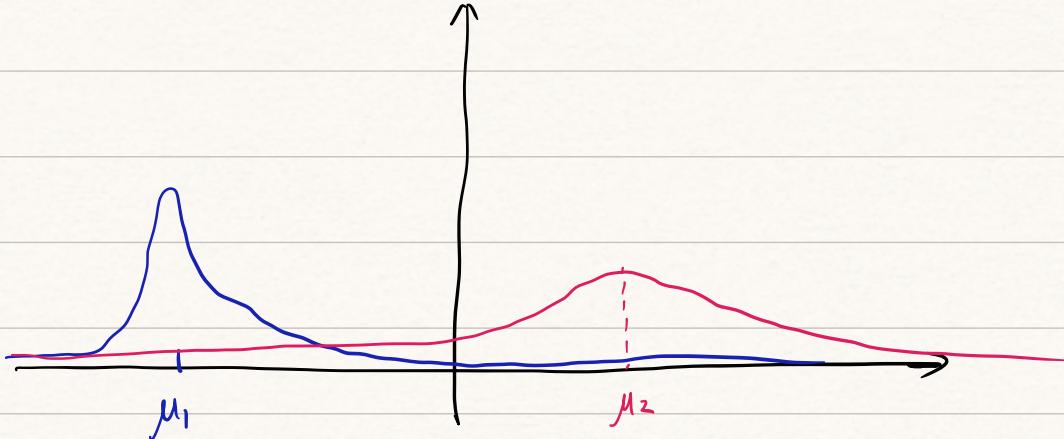
$$② N(x; \mu_j, \Sigma_j) = \frac{1}{(2\pi)^d |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) \right\}$$

$\{\mu_j\}_{j=1}^m$  : component mean

$\{\Sigma_j\}_{j=1}^m$  : component covariance matrix  $\Rightarrow$  PD matrix

$\Theta = \{P(j)\}_{j=1}^m, \{\mu_j\}_{j=1}^m, \{\Sigma_j\}_{j=1}^m\} \Rightarrow$  total # is  $3m$

$d=1$



Example:  $m^3$  different neighbourhoods  $\left\{ \begin{array}{l} \text{poor} \\ \text{middle-class} \\ \text{rich} \end{array} \right.$

$t=1, \dots, n$ ,  $n$ : # of houses

$x_t = \begin{bmatrix} x_{t1} \\ x_{t2} \end{bmatrix}$   $x_{t1}$ : price of  $t$ -th house  
 $x_{t2}$ : size of  $t$ -th house

Model the dataset  $\mathcal{D} = \{x_1, \dots, x_n\}$  as being drawn from GMM

Question:

- How many types of neighbourhoods? Model Selection
- Given a particular nbh.  $X$ , what is the  $\left\{ \begin{array}{l} \text{mean} \\ \text{Covariance} \end{array} \right.$   
 $\{(\mu_j, \Sigma_j)\}_{j=1}^n$
- In  $\mathcal{D}$ , what fraction of  $x_t$  comes from nbh. of a certain type?  $\{P(j)\}_{j=1}^m$

General Mixture Model: can not be GMM

$$P(X; \Omega) = \sum_{j=1}^m P(j) P(X | \theta_j)$$

pick parameter

→ task: Generating Models

Likelihood of observing  $\mathcal{D}$  given parameters  $\underline{\theta} = \{ \{P(j)\}, \{\underline{\theta}_j\} \}$

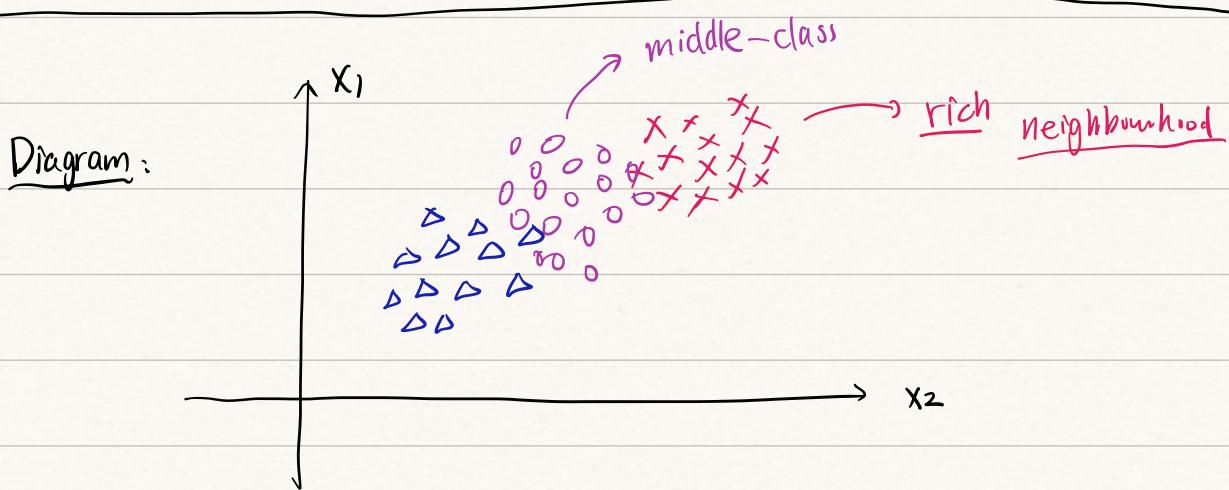
↓  
 $\underline{\theta} = \{x_1, \dots, x_n\}$

$$\mathcal{L}(\mathcal{D}; \theta) = \prod_{t=1}^n P(x_t; \theta)$$

$$\underline{\theta}_j = \{\mu_j, \Sigma_j\}$$

for GMM

$$= \prod_{t=1}^n \left( \sum_{j=1}^m P(j) P(x_t | \theta_j) \right)$$



Exercise:  $\mathcal{L}(\mathcal{D}; \theta) = \sum_{j=1}^m P(j) \left( \prod_{t=1}^n P(x_t | \theta_j) \right)$

↓  
 all the houses in certain type j

Story: All the houses in  $\mathcal{D}$  are of the same type,

but we DO NOT know which type of house represents  $\mathcal{D}$ .

Collaborative Filtering

→ Another Model Example

{ n users  
m movies

user	1	... M ...
	...	
item	n	... m ...
	...	

$r_{ij} \in \{1, \dots, 5\} \cup \{\text{Missing}\}$

→ sparse

## movie

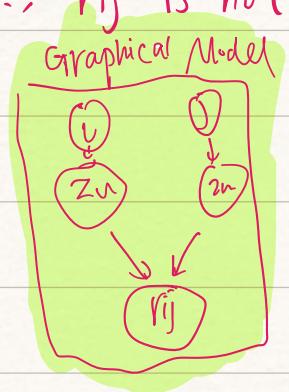
Each of the  $n$  users rates a subset of the  $m$  movies.

Provide  $r_{ij} \in \{1, \dots, 5\}$  or MISSING

### ★ defn of $I_D$

We say that  $(i, j) \in I_D$  if  $r_{ij}$  is observed, i.e.,  $r_{ij}$  is not missing. ( $\mathcal{D} \Rightarrow$  set of observed ratings)

$$\mathcal{D} = \{r_{ij} : \text{user } i \text{ rated movie } j\}$$



{ Types of movies indexed by  $z_m \in \{1, \dots, K_m\}$  }  $\Rightarrow$  to make problem simpler!  
{ Types of users indexed by  $z_u \in \{1, \dots, K_u\}$  }

↓  
degenerate!

How is a rating generated?

$$P(r_{ij} | i, j, \theta) = \sum_{z_u=1}^{K_u} \sum_{z_m=1}^{K_m} P(z_u | i) P(z_m | j) P(r_{ij} | z_u, z_m)$$

prob that movie  $j$   
 is of type  $z_m$   
 ↓  
 pnb that user  $i$   
 is of type  $z_u$   
 ↓  
 the rating given that  
 { user is of type  $z_u$   
 movie is of type  $z_m$

$$\underline{\theta} = \{ \{P(r|z_u, z_m)\}, \{P(z_u|i)\}, \{P(z_m|j)\} \}$$

Likelihood of a dataset  $\mathcal{D} = \{r_{ij} : \text{user } i \text{ rate movie } j\}$

$$L(\mathcal{D}; \underline{\theta}) = \prod_{(i,j) \in I_D} P(r_{ij} | i, j, \underline{\theta})$$

$$= \prod_{(i,j) \in I_D} \sum_{z_u} \sum_{z_m} P(z_u | i) P(z_m | j) P(r_{ij} | z_u, z_m)$$

How many Parameters  $\hat{\theta}$ ?

$$\textcircled{1} \quad \underbrace{\{P(z_u | i)\}}_{\substack{\downarrow \\ \text{mixture proportion}}} \quad \begin{cases} i=1, \dots, n \Rightarrow \text{user number} \\ z_u = 1, \dots, K_u \Rightarrow \text{user type number} \end{cases}$$

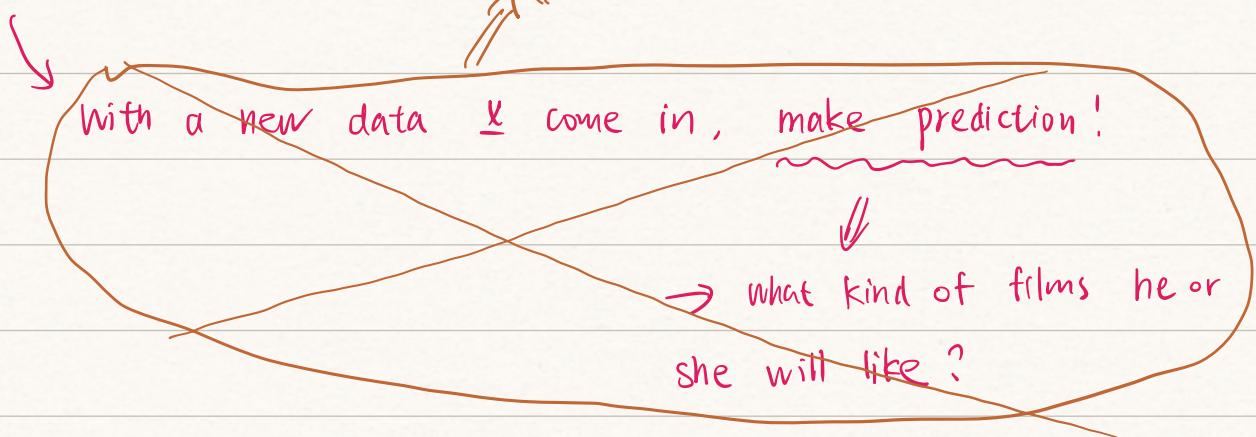
$\Rightarrow$  total # if  $n \cdot (K_u - 1)$

$$\textcircled{2} \quad \underbrace{\{P(z_m | j)\}}_{\substack{\uparrow \\ \rightarrow m \cdot (K_m - 1)}} \quad \xrightarrow{\text{自由参数}}$$

$$\textcircled{3} \quad \{P(r | z_u, z_m)\} \rightarrow (5-1) K_u K_m$$

What is our aim?

fill-in missing value instead!



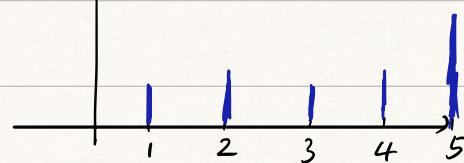
Once I have learnt  $\hat{\theta}$  from  $\mathcal{D} = \{r_{ij} : i \text{ rated movie } j\}$

Take user  $i^* = 1, \dots, n$  &  $r_{i^*j^*}$  is not available (missing)

$$\begin{array}{ccc} P(z_u | i^*) & P(z_m | j^*) & P(r_{ij} | z_u, z_m) \\ \uparrow & \uparrow & \\ \text{* 之前以前的投} & & \\ \text{票倾向} & & \end{array}$$

to predict  $P(r_{i^*j^*} | i^*, j^*)$

$$P(r_{i^*j^*} | i^*, j^*)$$



Cold-start problem: New user who has not rated any movies comes in.

→ using social media resources  
to attain more information

## Parameter Estimation

GMM:

$$P(x; \theta) = \sum_{j=1}^m P(j) N(x_j; \mu_j, \Sigma_j)$$

\* Likelihood  
注释: 这两个MLE是有明显不同的  
⇒ Complete Data 退化成了最简单的情形 (separable)

$$\mathcal{D} = \{x_t\}_{t=1}^n = X$$

→ Complete Dataset

① Pretend we know  $x_j$  belongs to  $\circ$  sub-mixture

Pretend that each observation  $x_1, \dots, x_n$  has additional information

about component responsible for generating it

$$j_1, \dots, j_n$$

→ Generating Information

$j_t \in \{1, \dots, m\}$ : component index of sample  $t$

$$\text{Def: } \delta(j|t) = \begin{cases} 1, & \text{if } j=j_t \\ 0, & \text{if } j \neq j_t \end{cases}$$

$$\mathcal{L}(\mathcal{D}; \theta) = \prod_{t=1}^n \sum_{j=1}^m P(j) N(x_t; \mu_j, \Sigma_j) \Rightarrow \text{difficult to deal with}$$

complete log-likelihood (Log-likelihood for complete data)

Complete data →  $(X, Z) = \{x_t, j_t\}_{t=1}^n$

which component generates  $x_t$ !

Very Easy to Deal with

$$l(\underbrace{x_1, \dots, x_n}_{X}, \underbrace{j_1, \dots, j_n}_{Z}; \theta) = \sum_{t=1}^n \log \left[ P(j_t) N(x_t; \mu_{j_t}, \Sigma_{j_t}) \right]$$

\* we assume that we have known  $x_t \leftarrow$  the  $j_t$ -th mixture

$$\begin{aligned}
 &= \sum_{t=1}^n \sum_{j=1}^m \delta(j|t) \log [P(j) N(x_j; \mu_j, \Sigma_j)] \\
 &\quad \text{筛选} \\
 &= \sum_{j=1}^m \sum_{t=1}^n \delta(j|t) \log P(j) \\
 &\quad + \sum_{j=1}^m \sum_{t=1}^n \delta(j|t) \log N(x_j; \mu_j, \Sigma_j)
 \end{aligned}$$

derivative  
 ↓  
 good property of  
 separation

Claim: Maximize  $l(X, Z; \theta)$  w.r.t  $\underline{\theta} = \{ \{P(j)\}, \{\mu_j\}, \{\Sigma_j\} \}$

yields that :

$$\hat{P}(j) = \frac{\hat{n}(j)}{n} \quad \hat{\mu}_j = \frac{1}{\hat{n}(j)} \sum_{t=1}^n \delta(j|t) x_t$$

$$\hat{\Sigma}_j = \frac{1}{\hat{n}(j)} \sum_{t=1}^n \delta(j|t) (x_t - \hat{\mu}_j)(x_t - \hat{\mu}_j)^T$$

Pf: (① and ②)

$$\textcircled{1} \text{ Note that } \sum_{j=1}^m P(j) = 1.$$

$x_1, \dots, x_n$  属于  $j$ -submodel 的数量.

$$\begin{aligned}
 L(\{P(j)\}, \lambda) &= \sum_{j=1}^m \sum_{t=1}^n \delta(j|t) \log P(j) + \lambda \left( 1 - \sum_{j=1}^m P(j) \right) \\
 &\quad \text{↓ } \hat{n}(j) \\
 &= \sum_{j=1}^m \hat{n}(j) \log P(j) + \lambda \left( 1 - \sum_{j=1}^m P(j) \right)
 \end{aligned}$$

$$\frac{\partial L}{\partial P(j)} = 0 \Rightarrow \frac{\hat{n}(j)}{P(j)} = \lambda \Rightarrow P(j) = \frac{\hat{n}(j)}{\lambda} \Rightarrow \text{solve for } \lambda = n$$

$\Rightarrow$  optimal  $P(j)$  must be proportional to  $\hat{n}(j)$

$$\Rightarrow \hat{P}(j) = \frac{\hat{n}(j)}{n}$$

$$\textcircled{2} \quad \ell(X, Z; \underline{\theta}) \stackrel{c}{=} \sum_{j=1}^m \sum_{t=1}^n \delta(j|t) \log [N(x_t; \mu_j, \Sigma_j)]$$

$$\stackrel{c}{=} \sum_{j=1}^m \sum_{t=1}^n \delta(j|t) \left[ -\frac{1}{2} (x_t - \mu_j)^T \Sigma_j^{-1} (x_t - \mu_j) \right]$$

$$\text{Then } \frac{\partial \ell}{\partial \mu_j} = 0 \Rightarrow \sum_{t=1}^n \delta(j|t) [\Sigma_j^{-1} (x_t - \mu_j)] = 0$$

$$\Rightarrow \sum_{t=1}^n \delta(j|t) (x_t - \mu_j) = 0$$

$$\Rightarrow \boxed{\hat{\mu}_j = \frac{1}{\hat{n}(j)} \sum_{t=1}^n \delta(j|t) x_t}$$

↓ just calculate in the neighbourhood!

while doing ③

we must notice that  $\Sigma_j \geq 0$

What if we cannot have the Ground-truth label?



we only have  $X = \{x_t\}_{t=1}^n \leftrightarrow \underbrace{\text{we don't have } \delta(j|t)}$

Method:

Replace them with estimation!  $(\underline{\theta}^{(e)})$



use  $\underline{\theta}^{(e)}$  to estimate  $\delta(j|t)$

$$\underbrace{P(j|x_t, \underline{\theta}^{(e)})}_{\delta(j|t)} = \frac{P(j, x_t | \underline{\theta}^{(e)})}{P(x_t | \underline{\theta}^{(e)})} \quad \text{use } (\mu_j^{(e)}, \Sigma_j^{(e)})$$

$$= \frac{p^{(e)}(j) P(x_t | j, \underline{\theta}^{(e)})}{\sum_{j=1}^m p^{(e)}(j) P(x_t | j, \underline{\theta}^{(e)})}$$

\* → Highlight

Instead of the HARD Assignments  $\{\delta(j|t)\}$ , we use the soft Assignments  $\{p^{(e)}(j|t)\}$  in  $e$ -th iteration

→ EM Alg.

Step 0: Initialize  $\underline{\theta}^{(0)}$

$$\text{e.g. } P(j) = \frac{1}{m} \text{ for } \forall j$$

$\mu_j$ : randomly chosen

$\Sigma_j$ : randomly chosen PD.

Step  $e$ . (E-step)

- Evaluate the Posterior assignment probability

$$P^{(e)}(j|t)$$

→ substitute for  $\delta(j|t)$

$$\delta(j|t)$$

Step  $e$  (M-step)

$$\hat{n}(j) = \sum_{t=1}^n P^{(e)}(j|t)$$

↓ hidden variable

- Update the parameters  $\underline{\theta}^{(e+1)}$ .

$$\hat{P}^{(e+1)}(j) = \frac{\hat{n}(j)}{n} \quad \hat{\mu}_j^{(e+1)} = \frac{1}{\hat{n}(j)} \sum_{t=1}^n P^{(e)}(j|t) x_t$$

$$③ \sum_j^{(l+1)} = \frac{1}{n(y_j)} \sum_{t=1}^n p^{(l)}(j|t) (\mathbf{x}_t - \hat{\mu}_j^{(l+1)}) (\mathbf{x}_t - \hat{\mu}_j^{(l+1)})^\top$$

i.e., Replaced  $\delta(j|t)$  with soft assignment  $p^{(l)}(j|t)$

Thm: (Main Thm)

Let  $\ell(\mathbf{x}; \theta^{(l)}) = \sum_{t=1}^n \log P(\mathbf{x}_t; \theta^{(l)})$ , the incomplete likelihood.



$$\ell(\mathbf{x}; \theta^{(l+1)}) \geq \ell(\mathbf{x}; \theta^{(l)}) \quad \forall l \in \mathbb{N}_0$$



Non-decreasing