

MODEL-BASED CLUSTERING

- SAME SETTING AND GOAL AS CENTER-BASED CLUSTERING
- STATISTICAL MODEL OF CLUSTERING

GIVEN A SEQUENCE $f_1(x), \dots, f_k(x)$ ON PROBABILITY

DENSITIES/DISTRIBUTION AT PROBABILITIES $p_1, \dots, p_k \geq 0$ $p_1 + \dots + p_k = 1$,

SUPPOSE $\vec{x}_1, \dots, \vec{x}_N$ WERE GENERATED IN THE

FOLLOWING

- ASSIGN CLUSTER LABEL, Z_i

SET $Z_i = z_i$ w/ PROBABILITY p_i

$$Z_i \in \{1, \dots, k\}$$

- DRAW \vec{x}_i CONDITIONAL ON $Z_i = z_i$

$$x_i \sim f_{z_i}(x)$$

- REPEAT INDEPENDENTLY FOR $i=1, \dots, N$

EQUIVALENTLY, WE COULD DRAW x_i FROM THE MIXTURE DISTRIBUTION

$$f_{\text{mix}}(x) = \sum_{j=1}^k p_j f_j(x)$$

Ex: GAUSSIAN MIXTURE MODEL

$$\begin{aligned} f_{\text{mix}}(x | p_1, \dots, p_k, \vec{\mu}_1, \dots, \vec{\mu}_k, \Sigma_1, \dots, \Sigma_k) \\ = \sum_{j=1}^k \frac{p_j}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu}_j)^T \Sigma_j^{-1} (\vec{x} - \vec{\mu}_j)\right) \end{aligned}$$

Ex: CAN USE MIXTURE OF ANY TYPE OF DISTRIBUTION OR COMBINATIONS!

t, t^* : GAUSSIAN, ... CAUCHY

MAIN IDEA

FOR FIXED K AND DENSITIES f_1, \dots, f_K DEFINED IN TERMS OF PARAMETERS ϕ

- 1) FIND THE MLE, $\hat{\phi}$, OF ϕ TO IDENTIFY THE BEST MODEL FOR MIXTURE DENSITY (EXPECTATION-MAXIMIZATION)
- 2) USE $f_{\text{MIX}}(x | \hat{\phi})$ TO ASSIGN CLUSTER LABELS

CAN REPORT FOR DIFFERENT CHOICES OF K (IF JUSTIFIED)
AND COMPARE RESULTS (AIC, BIC)

NOTE: THE CHOICE OF DENSITIES f_1, \dots, f_K AND SIMPLIFYING ASSUMPTIONS ON PARAMETERS CAN GREATLY INFLUENCE THE PERFORMANCE OF MODEL-BASED CLUSTERING. FOR THIS LESSON, WE WILL FOCUS ON MIXTURES OF GAUSSIANS (ALTHOUGH THIS CHOICE STRUGGLES W/ NON-CONVEX CLUSTERS).

$$\begin{aligned} \bullet \quad f_{\text{MIX}}(x | \phi) & \quad \phi = (p_1, \dots, p_K, \vec{\mu}_1, \dots, \vec{\mu}_K, \Sigma_1, \dots, \Sigma_K) \\ & = \sum_{j=1}^K p_j \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu}_j)^T \Sigma_j^{-1} (\vec{x} - \vec{\mu}_j)\right) \end{aligned}$$

• LIKELIHOOD OF DATA

$$\begin{aligned} L(\phi | \vec{x}_1, \dots, \vec{x}_N) & = \prod_{i=1}^N f(\vec{x}_i | \phi) \quad \text{WE DON'T OBSERVE } Z_i \\ & = \prod_{i=1}^N \left(\sum_{j=1}^K p_j \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (\vec{x}_i - \vec{\mu}_j)^T \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j)\right) \right) \end{aligned}$$

• LOG-LIKELIHOOD OF DATA

$$\log L(\phi | \vec{x}_1, \dots, \vec{x}_N) = \sum_{i=1}^N \log \left(\sum_{j=1}^K \dots \right)$$

CHOOSING K BASED ON AIC/BIC

FOR NOW, ASSUME WE CAN FIND $\hat{\phi} = (\hat{\mu}_1, \dots, \hat{\mu}_k, \hat{\sigma}_1, \dots, \hat{\sigma}_k)$

MAXIMIZING, i.e.

$$\begin{aligned}\hat{L} &= \max_{\phi} L(\phi | \vec{x}_1, \dots, \vec{x}_n) \\ \widehat{\log L} &= \max_{\phi} \log L(\phi | \vec{x}_1, \dots, \vec{x}_n) \\ &= \log L(\hat{\phi} | \vec{x}_1, \dots, \vec{x}_n)\end{aligned} \quad \log \hat{L} = \widehat{\log L}$$

IMPORTANTLY, \hat{L} GIVES AN ASSESSMENT OF HOW WELL THE MODEL FITS THE DATA,

- LARGER \hat{L} INDICATES A BETTER FIT
- TYPICALLY \hat{L} CAN BE INCREASED BY ADDING MORE COMPONENTS

* INCREASING K RISKS OVERFITTING !!
(POOR GENERALIZABILITY)

* PENALIZE MODELS W/ TOO MANY COMPONENTS (PARAMETERS)

AKAIKE INFORMATION CRITERION (AIC)

TO COMPARE MODELS (DIFFERENT CHOICES OF K), WE CAN COMPUTE AIC SCORES

$$AIC = -2 \log \hat{L} + 2p, \quad p = \# \text{ PARAMETERS}$$

WHERE LOW AIC IS BETTER.

↑
PENALTY FOR OVERPARAMETERIZATION

OTHERS DEFINE

$$AIC = 2 \log \hat{L} - 2p, \quad \text{SO HIGH AIC IS BETTER}$$

BAYES INFORMATION CRITERION

SIMILAR TO AIC, BUT ALSO INCORPORATES SAMPLE SIZE

$$BIC = -2 \log \hat{L} + p \log N, \quad p = \# \text{ PARAMETERS}$$

AGAIN, SMALL BIC IS BETTER

$$\left(\begin{array}{l} \text{IN MELUST PACKAGE } BIC = 2 \log \hat{L} - p \log N \\ \text{SO WE WANT MAXIMAL BIC} \end{array} \right)$$

NOTE: BOTH AIC & BIC ARE RELATED TO ASYMPTOTIC NORMALITY

OF MLE $\hat{\phi}$ AS $N \rightarrow \infty$

- TYPICALLY NEED N^d SAMPLES
- OTHER ASSUMPTIONS REQUIRED WHICH DO NOT HOLD W/O ADDITIONAL RESTRICTIONS ON ϕ
- RELIES ON AN APPROXIMATION TO BAYES FACTORS USING BIC

$$\frac{P(x_1, \dots, x_n \mid k \text{ CLUSTERS})}{P(x_1, \dots, x_n \mid k+1 \text{ CLUSTERS})} = \frac{P(k \text{ CLUSTERS} \mid x_1, \dots, x_n) P(k+1 \text{ CLUSTERS})}{P(k+1 \text{ CLUSTERS} \mid x_1, \dots, x_n) P(k \text{ CLUSTERS})}$$

REQUIRES SPECIFICATION
ON ALL PARAMETERS (WITH SUITABLE PRIORS)

Ex: GAUSSIAN MIXTURE IN \mathbb{R}^p w/ K COMPONENTS

- $(K-1)$ PARAMETERS FOR p_1, \dots, p_K
- Kp PARAMETERS FOR $\vec{\mu}_1, \dots, \vec{\mu}_K$
- $\frac{Kp(p+1)}{2}$ PARAMS. FOR $\Sigma_1, \dots, \Sigma_K$ IF NO RESTRICTIONS ON $\Sigma_1, \dots, \Sigma_K$
- 1 PARAMETER FOR $\Sigma_1, \dots, \Sigma_K$ IF WE ASSUME $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \sigma^2 I$

Ex: NCI DATA ON CANCER > FILES > NUMERICAL EXAMPLE > CLUSTERING

Assigning Cluster Labels

Once $f_{\text{max}}(x|\phi)$ has been chosen, assign x_i to cluster j for which

$$\hat{p}_j f_j(x_i | \hat{\mu}_j, \hat{\Sigma}_j)$$

is maximized, i.e.

$$z_i = \underset{j=1, \dots, K}{\operatorname{argmax}} p_j f_j(x_i | \hat{\mu}_j, \hat{\Sigma}_j)$$

Finding \hat{L} : The Expectation Maximization Algorithm

If we knew original cluster labels z_1, \dots, z_n then

$$\begin{aligned} L(\phi | x_1, \dots, x_n, z_1, \dots, z_n) &= \prod_{i=1}^N \sum_{j=1}^K \mathbb{1}(z_i=j) f_j(\vec{x}_i | \hat{\mu}_j, \hat{\Sigma}_j) \\ &= \prod_{i=1}^N \prod_{j=1}^K \left[p_j f_j(x_i, \mu_j, \Sigma_j) \right]^{\mathbb{1}(z_i=j)} \\ &= \exp \left(\sum_{i=1}^N \sum_{j=1}^K \mathbb{1}(z_i=j) \left[\log(p_j) - \frac{1}{2} \log |\Sigma_j| - \frac{d}{2} \log(2\pi) \right. \right. \\ &\quad \left. \left. - \frac{1}{2} (\vec{x}_i - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (\vec{x}_i - \hat{\mu}_j) \right] \right) \end{aligned}$$

Idea behind EM is to

E-step: estimate $P(z_i=j | \vec{x}_1, \dots, \vec{x}_n, \phi)$ and construct $E_{z|x,\phi} [\log L(\phi | x_1, \dots, x_n, z_1, \dots, z_n)]$

M-step: maximize $E_{z|x,\phi} (\log(\phi | x, z))$

Then repeat to convergence (to a local optima)

Ex: Mixture of Two Gaussians

$$\phi = (p_1, p_2, \mu_1, \Sigma_1, \mu_2, \Sigma_2)$$

$$\log L(\phi | x_1, \dots, x_N, z_1, \dots, z_N) = \sum_{i=1}^N \sum_{j=1}^2 \mathbb{1}(z_i=j) \left[\log p_j - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\vec{x}_i - \vec{\mu}_j)^T \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j) - \frac{d}{2} \log(2\pi) \right]$$

1) INITIALIZE VALUES FOR ϕ ($p_1, p_2 \neq 0$)

2) E-STEP: LET $\phi^{(t)}$ BE OUR CURRENT ESTIMATE FOR ϕ . SET

$$T_{i,j}^{(t)} = P(z_i=j | \vec{x}_i, \phi^{(t)}) = \frac{p_j^{(t)} f_j(\vec{x}_i, \mu_j^{(t)}, \Sigma_j^{(t)})}{\sum_{j=1}^2 p_j^{(t)} f_j(\vec{x}_i, \mu_j^{(t)}, \Sigma_j^{(t)})} \propto p_j^{(t)} f_j(\vec{x}_i, \mu_j^{(t)}, \Sigma_j^{(t)})$$

AND LET

$$\begin{aligned} E_{z|x_1, \dots, x_N, \phi^{(t)}} [\log(\phi | x_1, \dots, x_N, z_1, \dots, z_N)] \\ = \sum_{i=1}^N \sum_{j=1}^2 T_{i,j}^{(t)} \left[\log p_j - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\vec{x}_i - \vec{\mu}_j)^T \Sigma_j^{-1} (\vec{x}_i - \vec{\mu}_j) - \frac{d}{2} \log(2\pi) \right] \\ = Q(\phi | \phi^{(t)}) \end{aligned}$$

3) M-STEP: MAXIMIZE $Q(\phi | \phi^{(t)})$ OVER ALL ϕ . CLOSED FORM EXPRESSIONS IN THIS CASE

$$p_j^{(t+1)} = \frac{\sum_{i=1}^N T_{i,j}^{(t)}}{\sum_{i=1}^N \sum_{j=1}^2 T_{i,j}^{(t)}} \propto \sum_{i=1}^N T_{i,j}^{(t)}$$

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^N T_{i,j}^{(t)} \vec{x}_i}{\sum_{i=1}^N T_{i,j}^{(t)}}$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^N T_{i,j}^{(t)} (\vec{x}_i - \hat{\mu}_j^{(t+1)}) (\vec{x}_i - \hat{\mu}_j^{(t+1)})^T}{\sum_{i=1}^N T_{i,j}^{(t)}}$$

4) REPEAT 2-3 UNTIL CONVERGENCE TO A LOCAL OPTIMUM IS REACHED.

NOTE: RESTRICTING $\Sigma_1 = \Sigma_2 = \sigma^2 I$ AND SETTING $\sigma^2 \rightarrow 0$ (KNOWN) HAS INTERESTING CONNECTIONS TO K-MEANS CLUSTERING. SEE PSET #7 PROBLEM 4.