**Case Study 2: Document Retrieval**

# Review:
# Mixtures of Gaussians

Machine Learning for Big Data
CSE547/STAT548, University of Washington
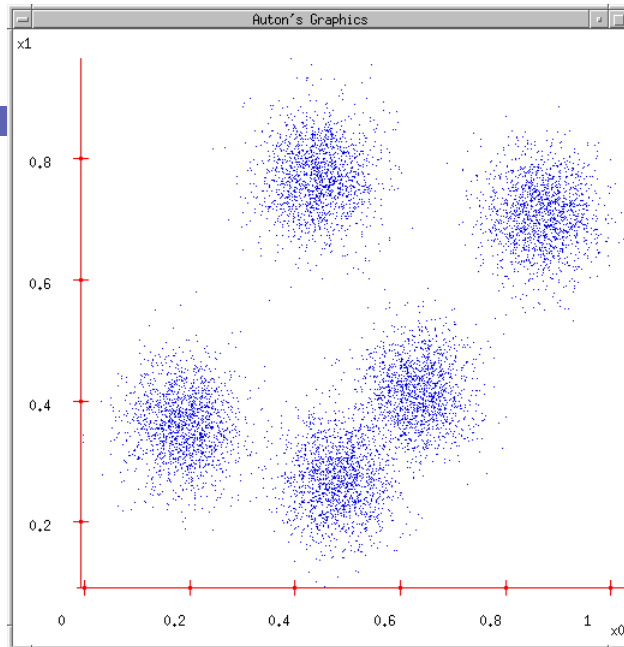
Emily Fox

January 28th, 2014

---

# Some Data
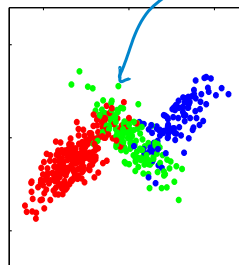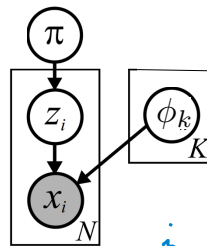
want to cluster

- unsupervised

- generative approach

# Gaussian Mixture Model

- Most commonly used mixture model
- Observations: $x^1, \cdots, x^N$   $x^i \in \mathbb{R}^d$

- Parameters:
  mix weights
  $$\pi = [\pi_1, \ldots, \pi_K]$$   $K$ # of clusters
  $$\phi = \{\phi_k\} = \{\mu_k, \Sigma_k\}$$
  params for cluster $k$
- Cluster indicator:
  $$z^i \in \{1, \ldots, K\} \qquad Pr(z^i = k) = \pi_k$$

- Per-cluster likelihood:
  $$N(x^i \mid \mu_k, \Sigma_k, z^i = k)$$



$x^i$

- Ex. $z^i$ = country of origin, $x^i$ = height of $i^{th}$ person
  - $k^{th}$ mixture component = distribution of heights in country $k$
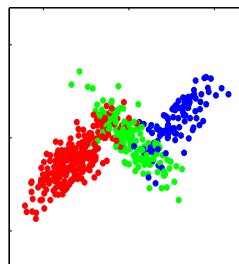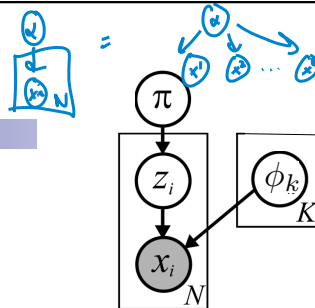
©Emily Fox 2014    3


# Generative Model

- We can think of *sampling* observations from the model

- For each observation $i$,
  - Sample a cluster assignment
    $$z^i \sim \pi \quad \text{"drawn from"}$$
    $1\ 2\ \cdots\ K$
  - Sample the observation from the selected Gaussian
    $$x^i \mid z^i \sim N(x^i \mid \mu_{z^i}, \Sigma_{z^i})$$

    can "generate" obs.



©Emily Fox 2014    4


2

# Also Useful for Density Estimation

$x^i \in \mathbb{R}^2$

$P(x^i | \theta)$

*Contour Plot of Joint Density*

bird's eye view
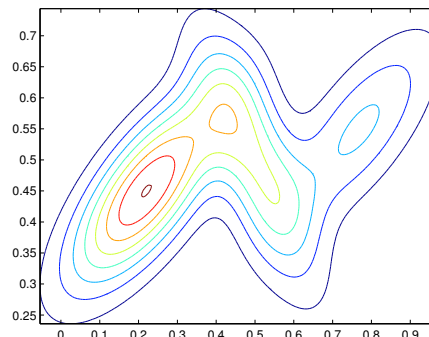
5

---

# Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

$k=$

*Mixture of 3 Gaussians*

*Contour Plot of Joint Density*

Each Gauss. has weight $\pi_k$ $\left( \sum \pi_k = 1 \right)$
and shape params $\{ \mu_k, \Sigma_k \}$

6

3

# Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians  $\{\pi_k\}$
  
  $(\pi_1, \dots, \pi_K)$  $\{\mu_k, \Sigma_k\}$

  *Mixture of 3 Gaussians*



$$p(x^i \mid \underbrace{\pi, \mu, \Sigma}_{\theta}) =$$

$$\sum_{k=1}^{K} \pi_k \, N(x^i \mid \mu_k, \Sigma_k)$$

$$\underset{P(z^i = k)}{\uparrow} \qquad \underset{P(x^i \mid z^i = k, \theta)}{\uparrow}$$

In 1D:

7

---

# Summary of GMM Components

- Observations $\qquad\qquad\qquad x_i \in \mathbb{R}^d, \quad i = 1, 2, \dots, N$

- Hidden cluster labels $\quad z_i \in \{1, 2, \dots, K\}, \quad i = 1, 2, \dots, N$

- Hidden mixture means $\qquad\qquad \mu_k \in \mathbb{R}^d, \quad k = 1, 2, \dots, K$

- Hidden mixture covariances $\quad \Sigma_k \in \mathbb{R}^{d \times d}, \quad k = 1, 2, \dots, K$

- Hidden mixture probabilities $\qquad\qquad \pi_k, \quad \sum_{k=1}^{K} \pi_k = 1$

***Gaussian mixture marginal and conditional likelihood* :**

$$p(x_i \mid \pi, \mu, \Sigma) = \sum_{z_i = 1}^{K} \pi_{z_i} \mathcal{N}(x_i \mid \mu_{z_i}, \Sigma_{z_i})$$

$$p(x_i \mid z_i, \pi, \mu, \Sigma) = \mathcal{N}(x_i \mid \mu_{z_i}, \Sigma_{z_i})$$

8

---

4

**Case Study 2: Document Retrieval**

# Application to Document Modeling

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

January 28th, 2014

**9**

---

# Task 2: Cluster Documents

- **Now:**
  - Cluster documents based on topic

"sports"

"world news"

**10**

# Document Representation

- Bag of words model

document *d*

previously:

$$x^d = \begin{bmatrix} \\ \\ \end{bmatrix}$$ ← vector fcn of word counts (e.g. tf-idf)

performed operations on this vector

now:

$$x^d = \{ w_1^d, \cdots, w_{N_d}^d \}$$ indices

unordered set of $N_d$ words with $w_i^d \in V$ vocab.
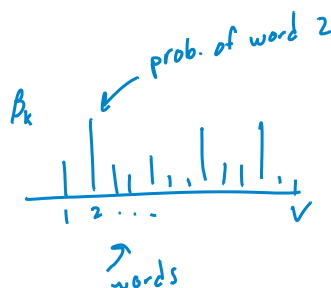
11

# A Generative Model

- Documents: $x^1, \cdots, x^D$ with $x^d = \{ w_1^d, \cdots, w_{N_d}^d \}$
- Associated topics: $z^1, \cdots, z^D$ with $z^d \in \{ 1, \cdots, K \}$
- Parameters: $\theta = \{ \pi, \beta \}$

$\uparrow$ # topics

as before $\begin{cases} \Pi = [ \pi_1, \cdots, \pi_K ] & \text{topic probabilities} \\ Pr(z^d = k) = \pi_k & \end{cases}$

$$\beta = \begin{array}{c} 1 \\ 2 \\ \vdots \\ K \end{array} \begin{bmatrix} 1 & 2 & \cdots & V \\ \hline & & \beta_1 & \\ & \vdots & & \\ & \beta_K & & \end{bmatrix}$$

$\beta_k$ ← prob. of word 2
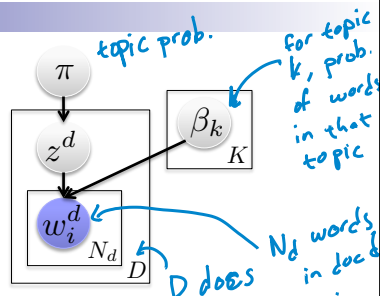
1 2 ... V

↗ words

12

6

# A Generative Model

- Documents: $x^1, \ldots, x^D$
- Associated topics: $z^1, \ldots, z^D$
- Parameters: $\theta = \{\pi, \beta\}$
- Generative model:

$z^d \sim \pi$    generate topic

$w_i^d \mid z^d \sim \beta_{z^d}$    $i = 1, \ldots, N_d$

Given topic $z^d = k$ for doc $d$, draw each word from $\beta_k$

$\pi$  topic prob.

for topic $k$, prob. of words in that topic

$z^d$    $\beta_k$    $K$

$w_i^d$    $N_d$    $D$

D docs

$N_d$ words in doc $d$

©Emily Fox 2014    13

---

# Form of Likelihood

- Conditioned on topic...

$$p(x^d \mid z^d, \beta) = \prod_{i=1}^{N_d} p(w_i^d \mid z^d, \beta) = \prod_{i=1}^{N_d} \beta_{z^d, w_i^d}$$

$\{w_1^d, \ldots, w_{N_d}^d\}$

- Marginalizing latent topic assignment:

$$p(x^d \mid \beta, \pi) = \sum_{k=1}^{K} \pi_k \, p(x^d \mid z^d = k, \beta_k)$$

$P(z^d = k)$

©Emily Fox 2014    14

# Review:
# EM Algorithm

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox
January 28th, 2014

©Emily Fox 2014                    **15**

---
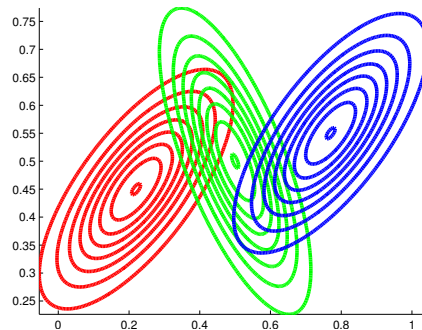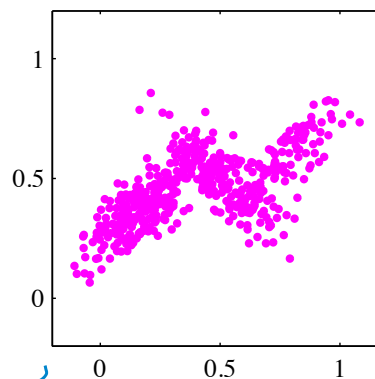
# Learning Model Parameters

- Want to learn model parameters

*Mixture of 3 Gaussians*                    *Our actual observations*

How???

from obs., estimate model params

*C. Bishop, Pattern Recognition & Machine Learning*

©Emily Fox 2014

# ML Estimate of Mixture Model Params

- Log likelihood

*assume $x^i$ iid*

$$L_x(\theta) \triangleq \log p(\{x^i\} \mid \theta) = \sum_i \log \sum_{z^i} p(x^i, z^i \mid \theta)$$

*← introduce cluster ind. + marg.*

- Want ML estimate

$$\hat{\theta}^{ML} = \arg\max_{\theta} L_x(\theta)$$

- Assume exponential family $p(x, z \mid \theta) = \dfrac{1}{Z(\theta)} e^{\theta' \phi(x,z)}$

$$L_x(\theta) = \sum_i \log \left( \sum_{z^i} e^{\theta^T \phi(z^i, x^i)} \right) - N \log Z(\theta)$$

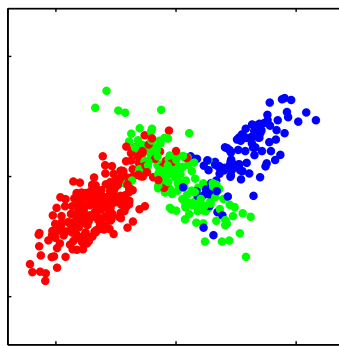- Neither convex nor concave and local optima

17

# Complete Data

- Imagine we have an assignment of each $x^i$ to a cluster

*Our actual observations*

*life would be easier...*

*Decouples into K ind. param. est. problems*



*Complete data labeled by true cluster assignments*

*"incomplete data"*

C. Bishop, Pattern Recognition & Machine Learning

9

# If "complete" data were observed…

- Assume class labels $z^i$ were observed in addition to $x^i$
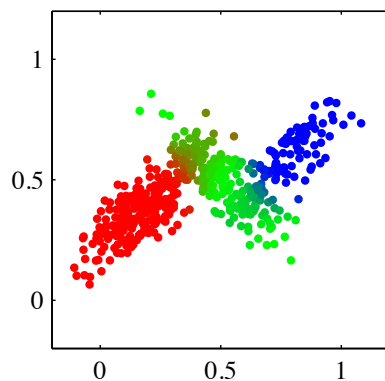
$$L_{x,z}(\theta) = \sum_i \log p(x^i, z^i \mid \theta)$$

- Compute ML estimates
  - Separates over clusters *k*!

- Example: mixture of Gaussians (MoG)   $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

---

# Cluster Responsibilities

- We must infer the cluster assignments from the observations

  *responsibility cluster k takes for obs i.*

  - Posterior probabilities of assignments to each cluster *given* model parameters:



$$r_{ik} = p(z^i = k \mid x^i, \pi, \phi) =$$

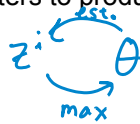$$= \frac{\pi_k \, p(x^i \mid \phi_k)}{\sum \pi_j \, p(x^i \mid \phi_j)}$$

*e.g.* $N(x^i \mid \mu_j, \Sigma_j)$

*Soft assignments to clusters*

*⭐ motivates iterative algorithm ⭐*

*C. Bishop, Pattern Recognition & Machine Learning*

# Iterative Algorithm

- Motivates a coordinate ascent-like algorithm:
  1. Infer missing values $z^i$ given estimate of parameters $\hat{\theta}$
  2. Optimize parameters to produce new $\hat{\theta}$ given "filled in" data $z^i$
  3. Repeat

  $z^i \circlearrowleft \theta$   max

- Example: MoG
  1. Infer "responsibilities"

  $$r_{ik}^{(t)} = p(z^i = k \mid x^i, \hat{\theta}^{(t-1)}) = \frac{\pi_k^{(t-1)}\, p(x^i \mid \phi_k^{(t-1)})}{\sum_j \pi_j^{(t-1)}\, p(x^i \mid \phi_j^{(t-1)})}$$

  2. Optimize parameters

  $$\max \text{ w.r.t. } \pi_k : \quad \pi_k^{(t)} = \frac{1}{N}\sum r_{ik}^{(t)} = \frac{r_k^{(t)}}{N} \leftarrow \text{soft counts!}$$

  $$\max \text{ w.r.t. } \phi_k :$$

  $$\mu_k^{(t)} = \frac{\sum r_{ik}^{(t)} x^i}{r_k^{(t)}} \leftarrow \text{weighted mean} \qquad \Sigma_k^{(t)} = \frac{1}{r_k^{(t)}}\sum r_{ik}^{(t)} x_i \cdot x_i^T - \mu_k^{(t)}\mu_k^{(t)T}$$
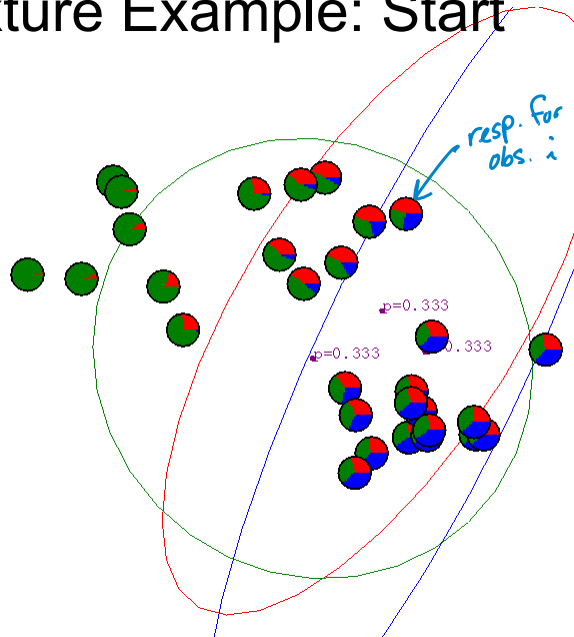
©Emily Fox 2014    21

---

# Gaussian Mixture Example: Start



resp. for obs. $i$

Initialize $\pi^{(0)}, \phi^{(0)}$

$\rightarrow$ compute $r_{ik}^{(1)}$

p=0.333   p=0.333   0.333

©Emily Fox 2014    22

# After first iteration

max like. given
soft counts

$\rightarrow \pi^{(1)}, \phi^{(1)}$

$\rightarrow$ new $r_{ik}^{(2)}$

p=0.306
p=0.273
=0.221

23

# After 2nd iteration

rinse +
repeat

p=0.374
p=0.306
=0.320

24

12

# After 3rd iteration



p=0.34

p=0.307

©Emily Fox 2014

25

# After 4th iteration



p=0.331

p=0.288

©Emily Fox 2014

26

13

# After 5th iteration

p=0.322
p=0.285

# After 6th iteration

p=0.315
p=0.287

# After 20th iteration

# Expectation Maximization (EM) – Setup

- More broadly applicable than just to mixture models considered so far

  *what we actually have*

- Model:   $x$   observable – *"incomplete" data*

  $y$   not (fully) observable – *"complete" data*     *what we wish we had*

  $\theta$   parameters

- Interested in maximizing (wrt $\theta$):

$$p(x \mid \theta) = \sum_y p(x, y \mid \theta)$$

  { *introduce complete data + marg.*

- Special case:

$$x = g(y)$$

  *e.g.*   $y = \begin{bmatrix} z \\ x \end{bmatrix}$ ← *class obs.*  ← *obs.*     *in standard mix model*

# EM Algorithm

- Initial guess: $\hat{\theta}^{(0)}$
- Estimate at iteration $t$: $\hat{\theta}^{(t)}$

- **E-Step**

  Compute $U(\theta, \hat{\theta}^{(t)}) = E\left[\log p(y|\theta) \mid x, \hat{\theta}^{(t)}\right]$

  (complete data)   (actual obs.)

- **M-Step**

  Compute $\hat{\theta}^{(t+1)} = \arg\max\limits_{\theta} U(\theta, \hat{\theta}^{(t)})$

  $$\Rightarrow L_x(\hat{\theta}^{(t+1)}) \geq L_x(\hat{\theta}^{(t)})$$

  $$\Rightarrow \hat{\theta}^{(t)} \text{ converges to a local mode}$$

---

# Example – Mixture Models

$E_{q_k}[I(z^i = k)]$
$= p(z^i = k \mid x^i, \hat{\theta}^{(t)})$
$\triangleq r_{ik}$

- **E-Step** Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y \mid \theta) \mid x, \hat{\theta}^{(t)}]$
- **M-Step** Compute $\hat{\theta}^{(t+1)} = \arg\max\limits_{\theta} U(\theta, \hat{\theta}^{(t)})$

- Consider $y^i = \{z^i, x^i\}$ i.i.d.

$$p(x^i, z^i \mid \theta) = \pi_{z^i} p(x^i \mid \phi_{z^i}) = \prod_{k=1}^{K} \left(\pi_k \, p(x^i \mid \phi_k)\right)^{I(z^i = k)}$$

$$E_{q_t}[\log p(y \mid \theta)] = \sum_i E_{q_t}[\log p(x^i, z^i \mid \theta)] =$$

$$= \sum_i \sum_k r_{ik} \log \pi_k + \sum_i \sum_k r_{ik} \log p(x^i \mid \phi_k)$$

E-step: computing $r_{ik}$'s

M-step: maximize wrt $\pi_k, \phi_k$

# Initialization

- In mixture model case where $y^i = \{z^i, x^i\}$ there are many ways to initialize the EM algorithm

- Examples:
  - ☐ Choose K observations at random to define each cluster. Assign other observations to the nearest "centriod" to form initial parameter estimates
  - ☐ Pick the centers sequentially to provide good coverage of data
  - ☐ Grow mixture model by splitting (and sometimes removing) clusters until K clusters are formed

- Can be quite important to convergence rates in practice

  *+ quality of local optima reached*

33

---

# What you need to know

- Mixture model formulation
  - ☐ Generative model
  - ☐ Likelihood
- Expectation Maximization (EM) Algorithm
  - ☐ Derivation ← *from prev. ML / Stat course*
  - ☐ Concept of non-decreasing log likelihood
  - ☐ Application to standard mixture models

34

---

17

**Case Study 2: Document Retrieval**

Review:
Connection to k-means

Machine Learning for Big Data
CSE547/STAT548, University of Washington
Emily Fox
January 28th, 2014

---

# K-means
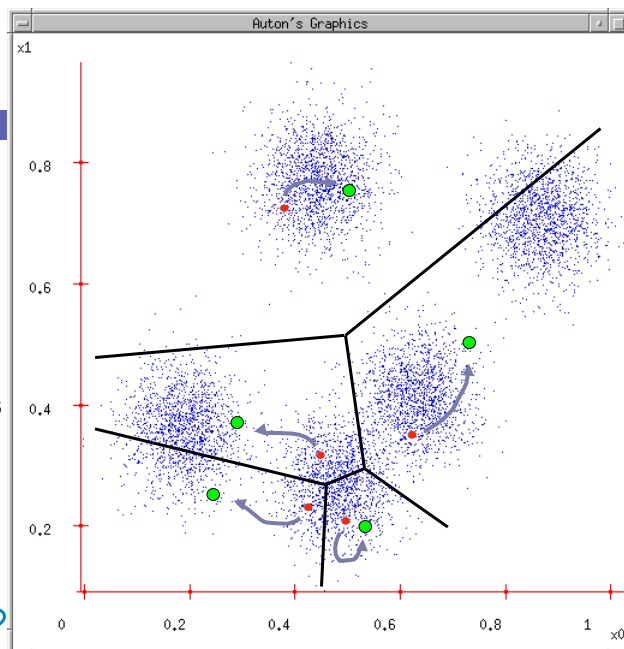*Recall*

1. Ask user how many clusters they'd like. *(e.g. k=5)*

2. Randomly guess k cluster Center locations

3. Each datapoint finds out which Center it's closest to.

4. Each Center finds the centroid of the points it owns

*iterative alg. making HARD assignments*

36

18

# K-means

- Randomly initialize $k$ centers
  - $\mu^{(0)} = \mu_1^{(0)}, \ldots, \mu_k^{(0)}$

- **Classify**: Assign each point $j \in \{1, \ldots m\}$ to nearest center:
  - $z^j \leftarrow \arg\min_i ||\mu_i - \mathbf{x}^j||_2^2$  ← hard assign.

- **Recenter**: $\mu_i$ becomes centroid of its point:
  - $\mu_i^{(t+1)} \leftarrow \arg\min_\mu \sum_{j:z^j=i} ||\mu - \mathbf{x}^j||_2^2$
  - Equivalent to $\mu_i \leftarrow$ average of its points!

37

# Special case: spherical Gaussians and hard assignments

$$P(z^i = k, \mathbf{x}^i|\theta) = \frac{1}{(2\pi)^{d/2} \| \Sigma_k \|^{1/2}} \exp\left[-\frac{1}{2}\left(\mathbf{x}^i - \mu_k\right)^T \Sigma_k^{-1}\left(\mathbf{x}^i - \mu_k\right)\right] P(z^i = k)$$

$\underbrace{\qquad\qquad}_{P(x^i \mid z^i = k, \theta)}$  $\underbrace{\qquad}_{P(z^i = k)}$

- If P(X|z=k) is spherical, with same $\sigma$ for all classes:

$$P(\mathbf{x}^i \mid z^i = k) \propto \exp\left[-\frac{1}{2\sigma^2}\left\|\mathbf{x}^i - \mu_k\right\|^2\right]$$

$P(x^i \mid z^i = k)$

assume $\Sigma_k = \begin{bmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{bmatrix} = \sigma^2 I$

- Then, compare EM objective with k-means:

EM: $\max_\theta \prod_i \sum_{z^i} P(x^i, z^i \mid \theta)$

maximizing marginal likelihood

k-means: $\max_{\{z^i\}, \theta} \prod_i P(x^i \mid z^i, \theta)$

OR if $\pi_k = \frac{1}{k}$  $\forall k$

$\max_{\{z^i\}, \theta} \prod_i P(x^i, z^i \mid \theta)$

$\max_\theta \prod_i \max_{z^i} P(x^i, z^i \mid \theta)$

38

19