

Case Study 2: Document Retrieval

Review: Mixtures of Gaussians

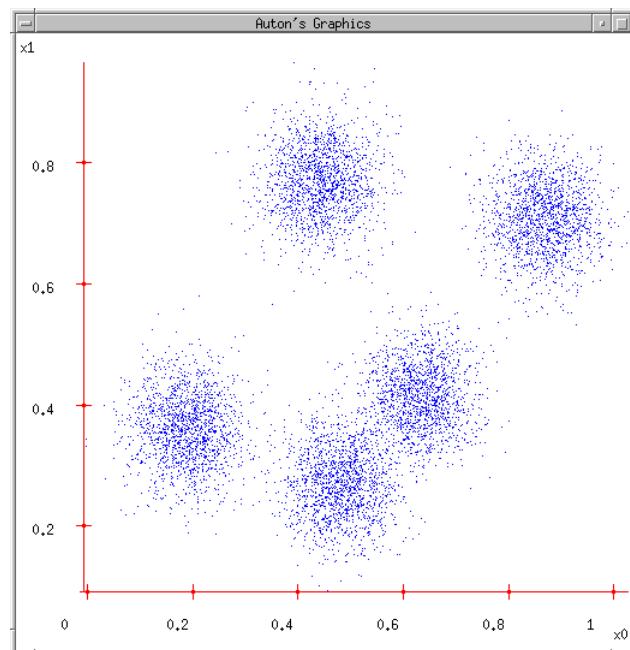
Machine Learning for Big Data
CSE547/STAT548, University of Washington
Emily Fox
April 21st, 2015

©Emily Fox 2015

1

Some Data

want to cluster
- unsup.
- generative approach

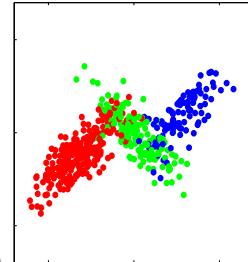
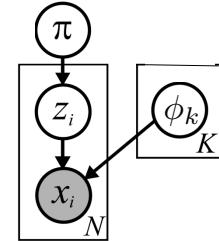


©Emily Fox 2015

2

Gaussian Mixture Model

- Most commonly used mixture model
- Observations: $x^1, \dots, x^N \quad x^i \in \mathbb{R}^d$
- Parameters: mix. weights
 $\pi = [\pi_1, \dots, \pi_K]$ # clusters
 $\phi = \{\phi_k\} = \{\mu_k, \Sigma_k\}$ params for cluster k
- Cluster indicator:
 $z^i \in \{1, \dots, K\} \quad \Pr(z^i = k) = \pi_k$
- Per-cluster likelihood:
 $N(x^i | \mu_k, \Sigma_k, z^i = k)$
- Ex. z^i = country of origin, x^i = height of i^{th} person
 - k^{th} mixture component = distribution of heights in country k

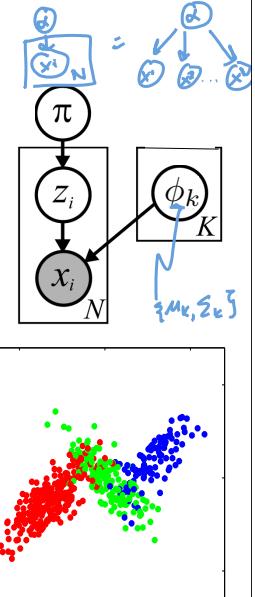


©Emily Fox 2015

3

Generative Model

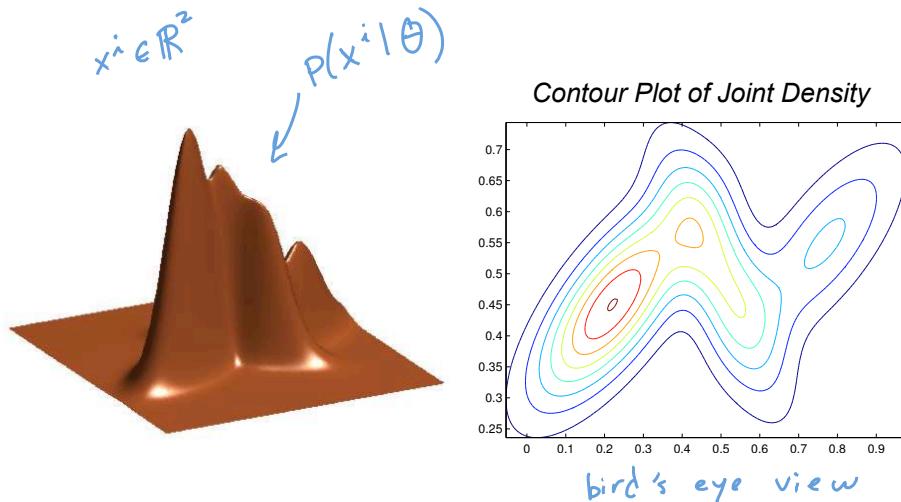
- We can think of sampling observations from the model
 - For each observation i ,
 - Sample a cluster assignment
 $z^i \sim \pi$ sample from or "drawn" from
 - Sample the observation from the selected Gaussian
 $x^i | z^i \sim N(x^i | \mu_{z^i}, \Sigma_{z^i})$
 ↑ given, "conditioned upon"
- can "generate" obs.



©Emily Fox 2015

4

Also Useful for Density Estimation

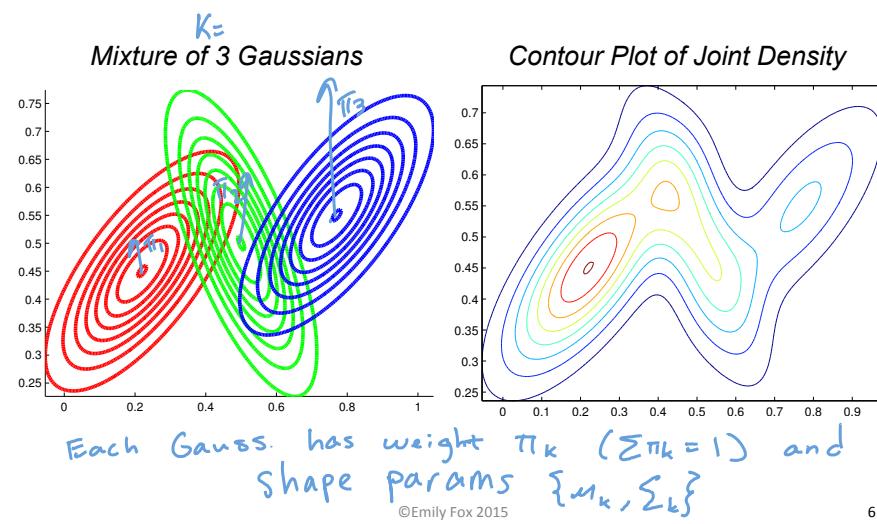


©Emily Fox 2015

5

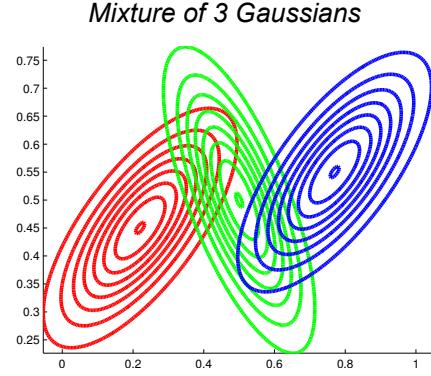
Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians



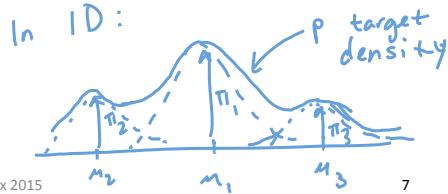
Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians



$$p(x^i | \pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k N(x^i | \mu_k, \Sigma_k)$$

↑ ↑
 $p(z^i | k)$ $p(x^i | z^i, \phi)$



©Emily Fox 2015

Summary of GMM Components

- Observations $x_i \in \mathbb{R}^d, i = 1, 2, \dots, N$
- Hidden cluster labels $z_i \in \{1, 2, \dots, K\}, i = 1, 2, \dots, N$
- Hidden mixture means $\mu_k \in \mathbb{R}^d, k = 1, 2, \dots, K$
- Hidden mixture covariances $\Sigma_k \in \mathbb{R}^{d \times d}, k = 1, 2, \dots, K$
- Hidden mixture probabilities $\pi_k, \sum_{k=1}^K \pi_k = 1$

Gaussian mixture marginal and conditional likelihood :

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

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

Case Study 2: Document Retrieval

Application to Document Modeling

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

Emily Fox

April 21st, 2015

©Emily Fox 2015

9

Task 2: Cluster Documents

Now:

- Cluster documents based on topic



©Emily Fox 2015

10

Document Representation

- Bag of words model

document d

previously
 $X = \begin{bmatrix} \end{bmatrix}$ ↗ vector func
 of word counts
 (e.g. tf-idf)
 performed operations on
 this vector

now :

$X = \{w_1, \dots, w_N\}$
 ↗ list
unordered set of N words in
 doc.
 $w_i \in V$ (vocab)

©Emily Fox 2015

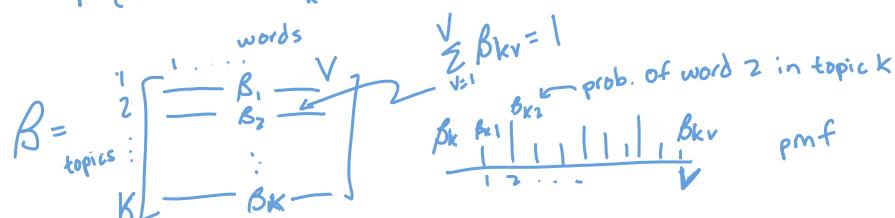
11

A Generative Model

- Documents: x^1, \dots, x^D ↗ # of docs with $x^d = \{w_1^d, \dots, w_{N_d}^d\}$ ↗ word ID (not count)
- Associated topics: z^1, \dots, z^D with $z^d \in \{1, \dots, K\}$ ↗ # of topics
- Parameters: $\theta = \{\pi, \beta\}$ ↗ mix weights ↗ cluster params ↗ # of words in doc d

$\pi = [\pi_1, \dots, \pi_K]$ topic prob.

$$\Pr(z^i = k) = \pi_k$$



©Emily Fox 2015

12

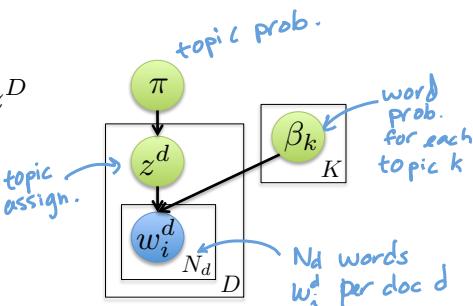
A Generative Model

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

Sample topic:
 $z^d \sim \pi$

Sample words:
 $w_i^d | z^d \sim \beta_{z^d} \quad i=1, \dots, N_d$

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



©Emily Fox 2015

13

Form of Likelihood

- Conditioned on topic... N_d iid words from topic z^d

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

- Marginalizing latent topic assignment:

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

©Emily Fox 2015

14

Case Study 2: Document Retrieval

Review: EM Algorithm

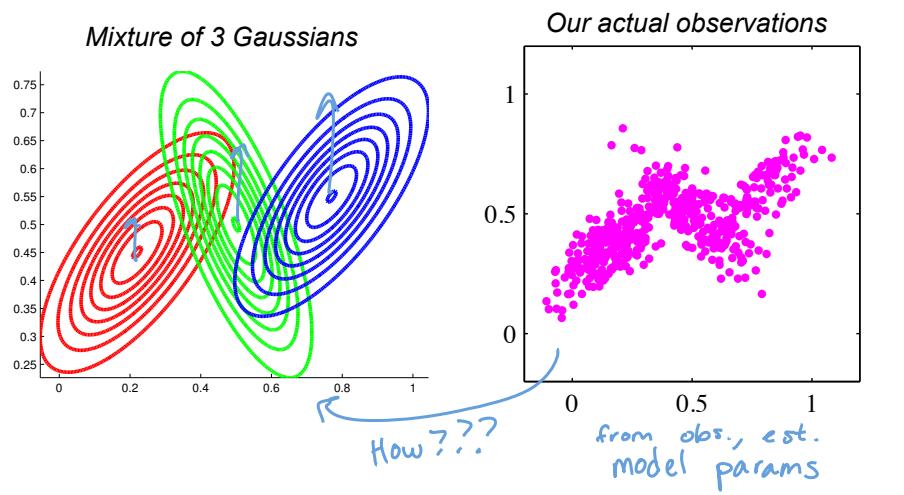
Machine Learning for Big Data
 CSE547/STAT548, University of Washington
 Emily Fox
 April 21st, 2015

©Emily Fox 2015

15

Learning Model Parameters

- Want to learn model parameters



C. Bishop, *Pattern Recognition & Machine Learning* 16

ML Estimate of Mixture Model Params

- Log likelihood

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

$p(x^i | \theta) = \prod_i p(x^i | \theta)$

- Want ML estimate

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

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

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

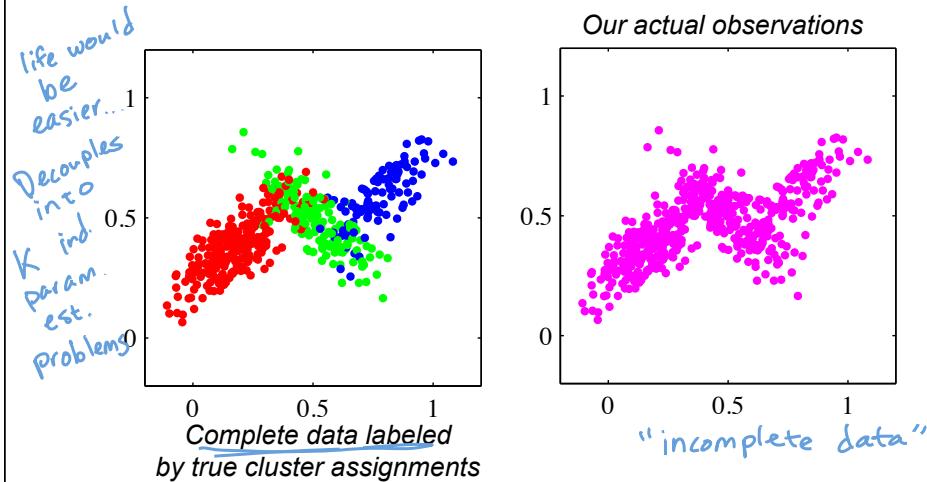
- Neither convex nor concave and local optima

©Emily Fox 2015

17

Complete Data

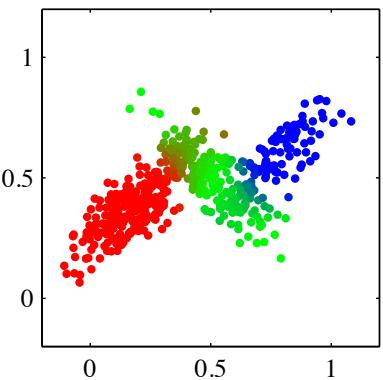
- Imagine we have an assignment of each x^i to a cluster



©Emily Fox 2015

C. Bishop, Pattern Recognition & Machine Learning 18

Cluster Responsibilities

- We must infer the cluster assignments from the observations
- 
- "responsibility" of cluster k for point i
- Posterior probabilities of assignments to each cluster *given* model parameters:
- $$r_{ik} = p(z^i = k | x^i, \pi, \phi) = \frac{\pi_k p(x^i | \phi_k)}{\sum_{j=1}^K \pi_j p(x^i | \phi_j)}$$
- e.g. $N(x^i | \mu_j, \Sigma_j)$
- Soft assignments to clusters
- * Motivates iterative alg *
- ©Emily Fox 2015 C. Bishop, Pattern Recognition & Machine Learning 19

Iterative Algorithm

- Motivates a coordinate ascent-like algorithm:
 - Infer missing values z^i given estimate of parameters $\hat{\theta}$
 - Optimize parameters to produce new $\hat{\theta}$ given "filled in" data z^i
 - Repeat
 - Example: MoG (derivation soon... + HW)
 - Infer "responsibilities"
$$r_{ik}^{(t)} = p(z^i = k | x^i, \hat{\theta}^{(t-1)}) = \frac{\pi_k^{(t-1)} p(x^i | \phi_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} p(x^i | \phi_j^{(t-1)})}$$
 - Optimize parameters

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

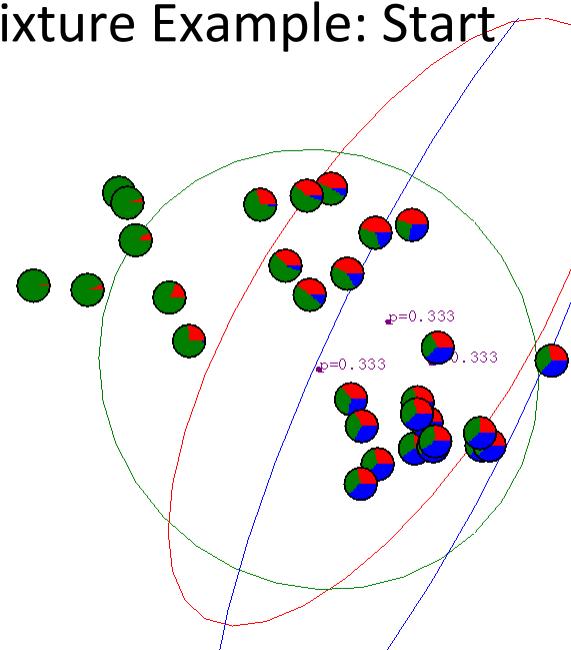
max w.r.t. ϕ_k :

$$\mu_k^{(t)} = \frac{\sum r_{ik}^{(t)} x^i}{r_k^{(t)}} \leftarrow \text{weighted mean}$$

$$\Sigma_k^{(t)} = \frac{1}{r_k^{(t)}} \sum r_{ik}^{(t)} x^i x^{i\top} - \frac{(\mu_k^{(t)}) (\mu_k^{(t)})^\top}{r_k^{(t)}}$$
- ©Emily Fox 2015

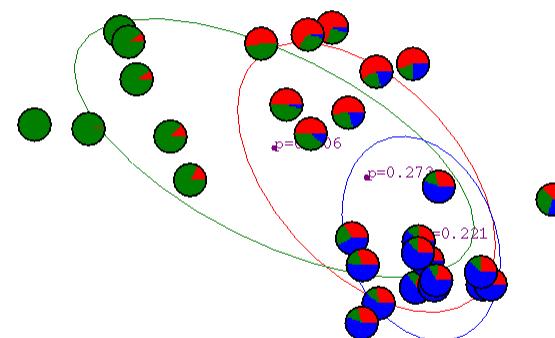
Gaussian Mixture Example: Start

Initialize
 $\pi^{(0)}, \phi^{(0)}$
 \rightarrow compute
 $r_{ik}^{(1)}$



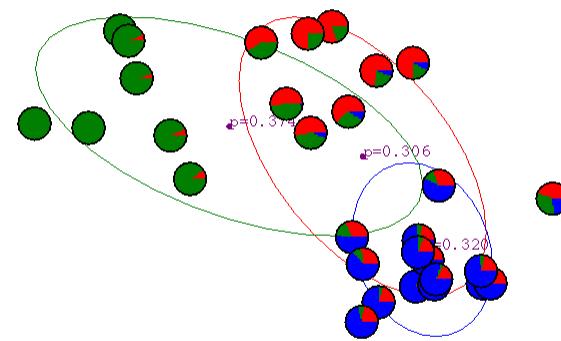
After first iteration

max like. given
soft counts
 $\rightarrow \pi^{(1)}, \phi^{(1)}$
 \rightarrow new $r_{ik}^{(2)}$



After 2nd iteration

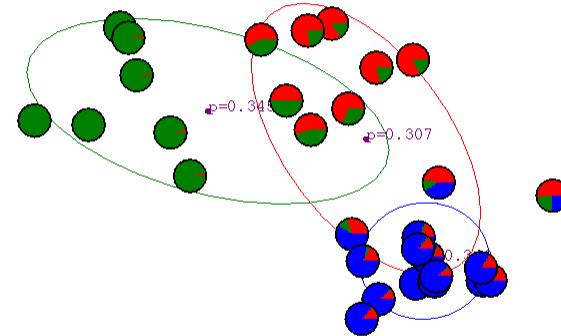
rinse +
repeat



©Emily Fox 2015

23

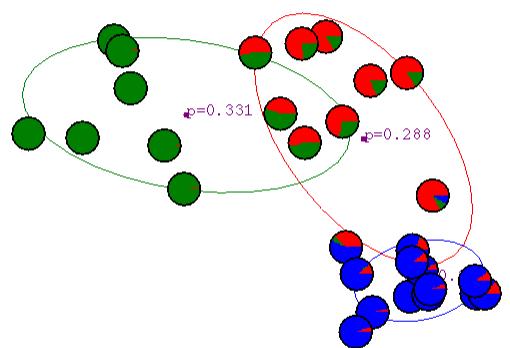
After 3rd iteration



©Emily Fox 2015

24

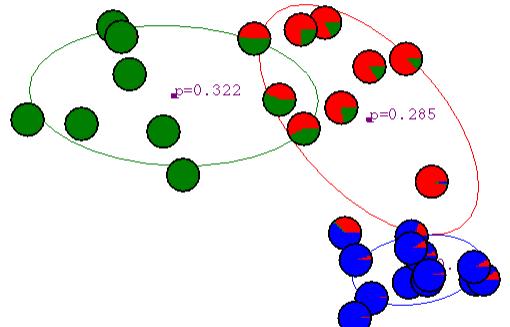
After 4th iteration



©Emily Fox 2015

25

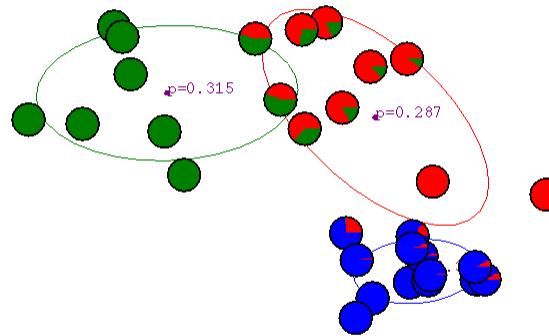
After 5th iteration



©Emily Fox 2015

26

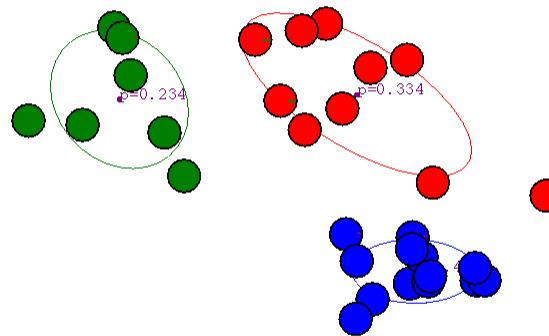
After 6th iteration



©Emily Fox 2015

27

After 20th iteration



©Emily Fox 2015

28

Expectation Maximization (EM) – Setup

- More broadly applicable than just to mixture models considered so far
- Model: x observable – “incomplete” data what we actually have
 y not (fully) observable – “complete” data ← what we wish we had
 θ parameters
- Interested in maximizing (wrt θ):

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

↑ introduce + marg. complete data

- Special case:

$$x = g(y)$$

deterministic fcn

e.g. $y = \begin{bmatrix} z \\ x \end{bmatrix}$

← class labels
← obs.

in standard mix model

©Emily Fox 2015

29

EM Algorithm

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

- E-Step**

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

- M-Step**

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

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

mild assumption $\hat{\theta}$ converges to a local mode

©Emily Fox 2015

30

Example – Mixture Models

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

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

$$\begin{aligned}
 p(x^i, z^i | \theta) &= \pi_{z^i} p(x^i | \phi_{z^i}) = \prod_{k=1}^K (\pi_k p(x^i | \phi_k))^{I(z^i=k)} \\
 E_{q_t} [\log p(y | \theta)] &= \sum_i E_{q_t} [\log p(x^i, z^i | \theta)] = \\
 &= \sum_{i \in K} r_{ik} \log \pi_k + \sum_{i \in K} \sum_{k \neq i} r_{ik} \log p(x^i | \phi_k) \\
 \text{M-step:} \quad \max_{\text{w.r.t. } \pi_k, \phi_k} &\quad \text{E-step:} \quad \text{compute the } r_{ik} \text{ based on } \hat{\theta}^{(t)} \\
 &\quad \left. \begin{aligned} E_{q_t} [I(z^i=k)] &= \\ &= p(z^i=k | x, \hat{\theta}^{(t)}) \\ &\approx r_{ik} \end{aligned} \right\} \text{fn of params}
 \end{aligned}$$

©Emily Fox 2015

31

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 mode

©Emily Fox 2015

32

What you need to know

- Mixture model formulation
 - Generative model
 - Likelihood
- Expectation Maximization (EM) Algorithm
 - Derivation
 - Concept of non-decreasing log likelihood
 - Application to standard mixture models

©Emily Fox 2015

33

Case Study 2: Document Retrieval

Review: Connection to k-means

Machine Learning for Big Data
CSE547/STAT548, University of Washington
Emily Fox
April 21st, 2015

©Emily Fox 2015

34

K-means

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*

K-means

- Randomly initialize k centers
 $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$
- **Classify:** Assign each point $j \in \{1, \dots, N\}$ to nearest center:

$$z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2 \quad \text{hard assign.}$$
- **Recenter:** μ_i becomes centroid of its point:

$$\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j:z^j=i} \|\mu - \mathbf{x}^j\|_2^2 \quad \text{update params.}$$
 - Equivalent to $\mu_i \leftarrow \text{average of its points!}$

©Emily Fox 2015 36

Special Case: Spherical Gaussians + hard assignments

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

- If $P(\mathbf{x}|z=k)$ is spherical, with same σ for all classes:

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

assume: $\Sigma_k = \begin{pmatrix} \sigma^2 & \dots & \sigma^2 \\ \vdots & \ddots & \vdots \\ \sigma^2 & \dots & \sigma^2 \end{pmatrix} = \sigma^2 I$

- Then, compare EM objective with k-means:

$\text{EM: } \max_{\theta} \prod_i \sum_{z^i} p(x^i, z^i \theta)$ <p style="margin-left: 100px;">maximizing marg. likelihood</p>	$\text{k-means: } \max_{\{\pi_k, \mu_k\}} \prod_i p(x^i z^i, \theta)$ <p style="margin-left: 100px;">OR if $\pi_k = \frac{1}{K} \forall k$</p> $\max_{\pi_k, \mu_k} \prod_i p(x^i, z^i \theta)$ $\max_{\theta} \prod_i \max_{z^i} p(x^i, z^i \theta)$
--	--

©Emily Fox 2015

37