

Case Study 2: Document Retrieval

Clustering Documents, Mixture Models, Expectation Maximization

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Emily Fox

January 29th, 2013

©Emily Fox 2013

1

Document Retrieval

- **Goal:** Retrieve documents of interest

- **Challenges:**

- Tons of articles out there
- How should we measure similarity?



©Emily Fox 2013

2

Task 1: Find Similar Documents

■ So far...

- **Input:** Query article
- **Output:** Set of k similar articles



©Emily Fox 2013

3

Task 2: Cluster Documents

■ Now:

- Cluster documents based on topic



©Emily Fox 2013

4

Document Representation

- Bag of words model



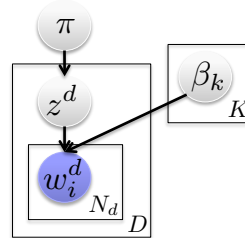
document d

A Generative Model

- Documents:
- Associated topics:
- Parameters: $\theta = \{\pi, \beta\}$

A Generative Model

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



©Emily Fox 2013

7

Form of Likelihood

- Conditioned on topic...

$$p(x^d | z^d, \beta) =$$

- Marginalizing latent topic assignment:

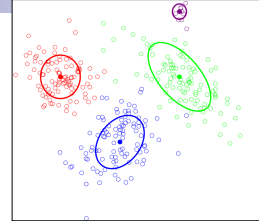
$$p(x^d | \beta, \pi) =$$

©Emily Fox 2013

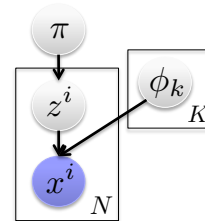
8

Gaussian Mixture Model

- Most commonly used mixture model
- Observations:
- Parameters:



- Likelihood:



- Ex. z^i = country of origin, x^i = height of i^{th} person
 - k^{th} mixture component = distribution of heights in country k

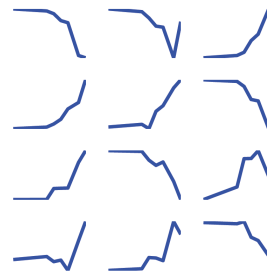
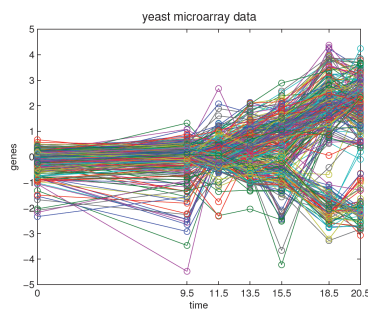
©Emily Fox 2013

9

Another Example

(Taken from Kevin Murphy's ML textbook)

- Data: gene expression levels
- Goal: cluster genes with similar expression trajectories



©Emily Fox 2013

10

Mixture models are useful for...

- Density estimation

- Allows for multimodal density

- Clustering

- Want membership information for each observation
 - e.g., topic of current document
- Soft clustering:

$$p(z^i = k | x^i, \theta) =$$

- Hard clustering:

$$z^{i*} = \arg \max_k p(z^i = k | x^i, \theta) =$$

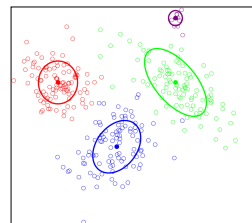
©Emily Fox 2013

11

Issues

- Label switching

- Color = label does not matter
- Can switch labels and likelihood is unchanged



- Log likelihood is not convex in the parameters

- Problem is simpler for “complete data likelihood”

©Emily Fox 2013

12

ML Estimate of Mixture Model Params

- Log likelihood

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

- Want ML estimate

$$\hat{\theta}^{ML} =$$

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

$$L_x(\theta) =$$

- Neither convex nor concave and local optima

©Emily Fox 2013

13

If “complete” data were observed...

- Assume class labels z^i were observed in addition to x^i

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

- Compute ML estimates
 - Separates over clusters k !

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

©Emily Fox 2013

14

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

- Example: MoG (derivation soon... + HW)

1. Infer “responsibilities”

$$r_{ik} = p(z^i = k \mid x^i, \hat{\theta}^{(t-1)}) =$$

2. Optimize parameters

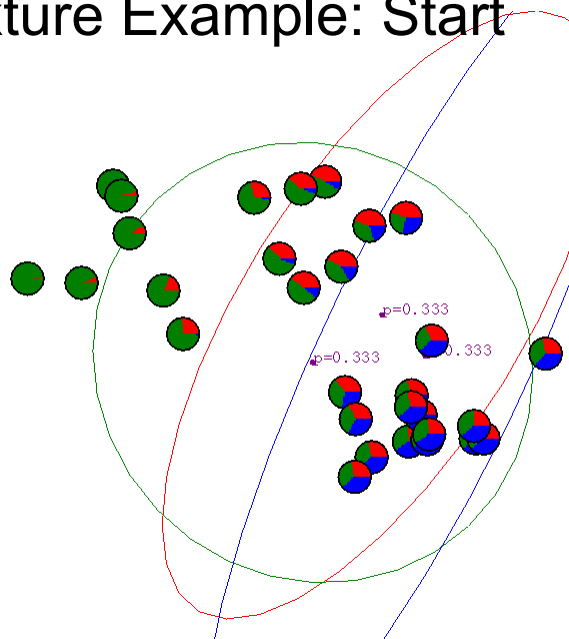
max w.r.t. π_k :

max w.r.t. ϕ_k :

©Emily Fox 2013

15

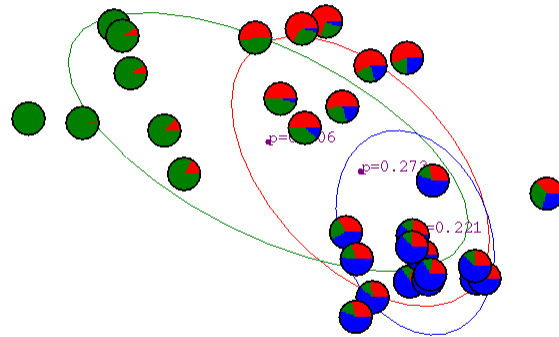
Gaussian Mixture Example: Start



©Emily Fox 2013

16

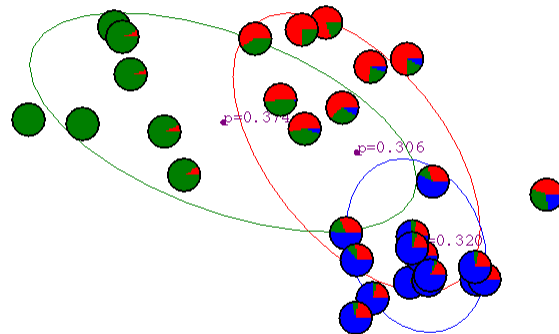
After first iteration



©Emily Fox 2013

17

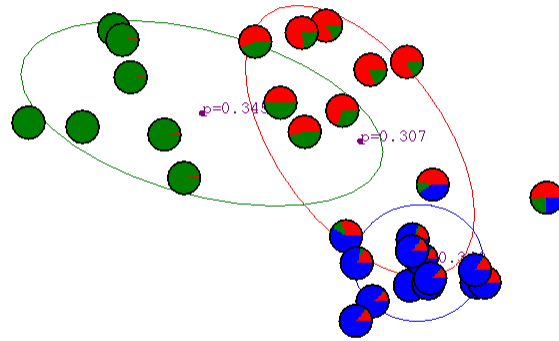
After 2nd iteration



©Emily Fox 2013

18

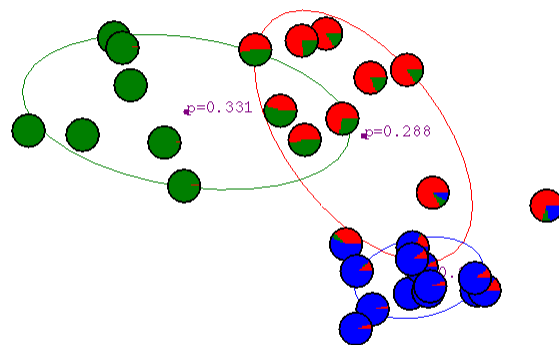
After 3rd iteration



©Emily Fox 2013

19

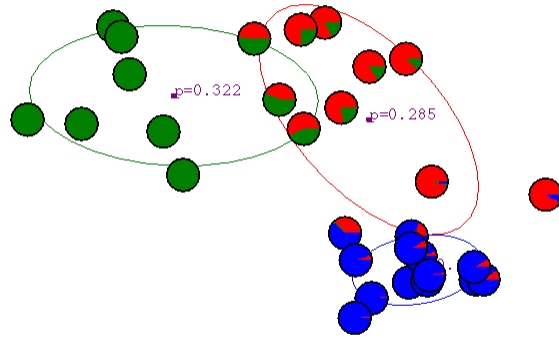
After 4th iteration



©Emily Fox 2013

20

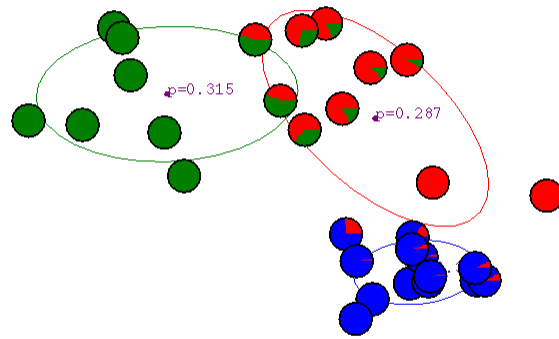
After 5th iteration



©Emily Fox 2013

21

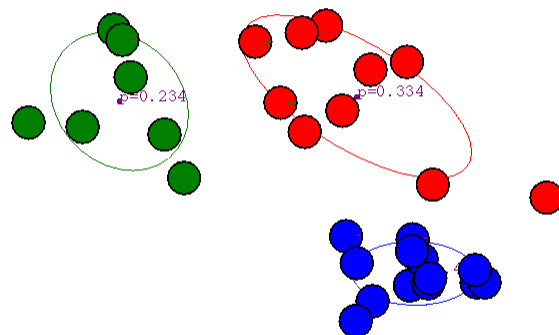
After 6th iteration



©Emily Fox 2013

22

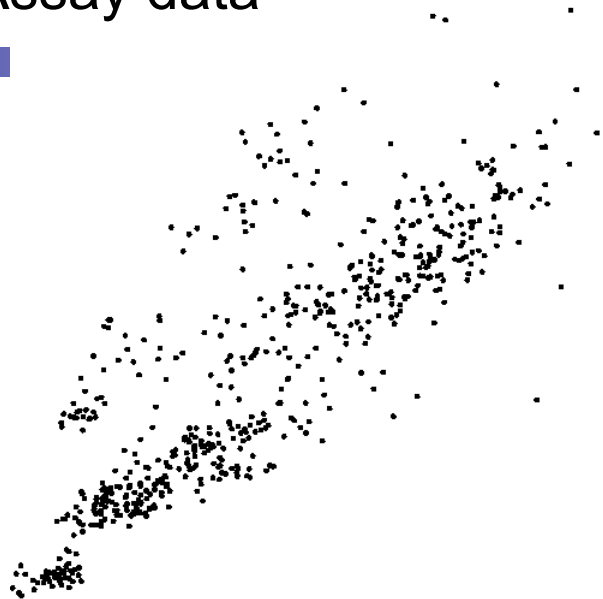
After 20th iteration



©Emily Fox 2013

23

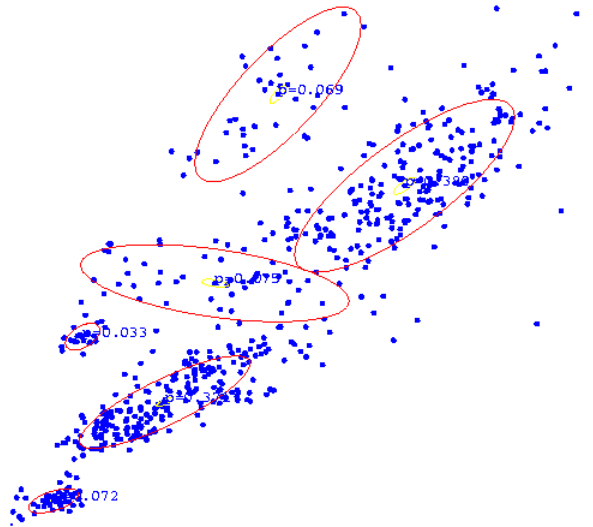
Some Bio Assay data



©Emily Fox 2013

24

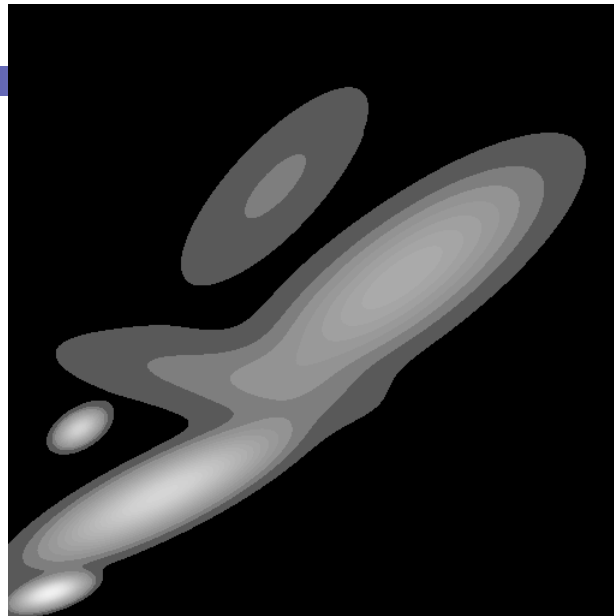
GMM clustering of the assay data



©Emily Fox 2013

25

Resulting Density Estimator



©Emily Fox 2013

26

Expectation Maximization (EM) – Setup

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

- Model: x observable – “incomplete” data
 y not (fully) observable – “complete” data
 θ parameters

- Interested in maximizing (wrt θ):

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

- Special case:

$$x = g(y)$$

©Emily Fox 2013

27

Expectation Maximization (EM) – Derivation

- Step 1
 - Rewrite desired likelihood in terms of complete data terms

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

- Step 2
 - Assume estimate of parameters $\hat{\theta}$
 - Take expectation with respect to $p(y | x, \hat{\theta})$

©Emily Fox 2013

28

Expectation Maximization (EM) – Derivation

- Step 3

- Consider log likelihood of data at any θ relative to log likelihood at $\hat{\theta}$

$$L_x(\theta) - L_x(\hat{\theta})$$

- **Aside: Gibbs Inequality** $E_p[\log p(x)] \geq E_p[\log q(x)]$

Proof:

Expectation Maximization (EM) – Derivation

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta, \hat{\theta}) - U(\hat{\theta}, \hat{\theta})] - [V(\theta, \hat{\theta}) - V(\hat{\theta}, \hat{\theta})]$$

- Step 4

- Determine conditions under which log likelihood at θ exceeds that at $\hat{\theta}$
Using Gibbs inequality:

If

Then

$$L_x(\theta) \geq L_x(\hat{\theta})$$

Motivates EM Algorithm

- Initial guess:
- Estimate at iteration t :

- **E-Step**

Compute

- **M-Step**

Compute

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.

$$p(x^i, z^i | \theta) = \pi_{z^i} p(x^i | \phi_{z^i}) =$$

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

Coordinate Ascent Behavior

- Bound log likelihood:

$$\begin{aligned} L_x(\theta) &= U(\theta, \hat{\theta}^{(t)}) + V(\theta, \hat{\theta}^{(t)}) \\ &\geq \\ L_x(\hat{\theta}^{(t)}) &= U(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) + V(\hat{\theta}^{(t)}, \hat{\theta}^{(t)}) \end{aligned}$$

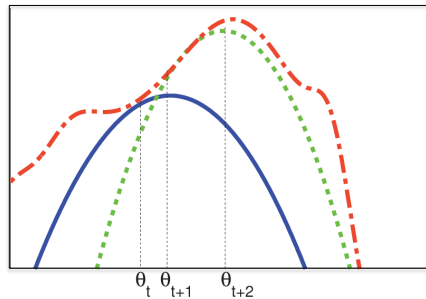


Figure from
KM textbook

©Emily Fox 2013

33

Comments on EM

- Since Gibbs inequality is satisfied with equality only if $p=q$, any step that changes θ should strictly **increase likelihood**
- In practice, can replace the **M-Step** with increasing U instead of maximizing it (**Generalized EM**)
- Under certain conditions (e.g., in exponential family), can show that EM **converges to a stationary point** of $L_x(\theta)$
- Often there is a **natural choice for y** ... has physical meaning
- If you want to choose any y , not necessarily $x=g(y)$, replace $p(y | \theta)$ in U with $p(y, x | \theta)$

©Emily Fox 2013

34

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 “centroid” 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

©Emily Fox 2013

35

MAP Estimation

- Bayesian approach:
 - Place **prior** $p(\theta)$ on parameters
 - Infer **posterior** $p(\theta | x)$
- Many, many, many motivations and implications
 - For the sake of this class, simplest motivation is to think of this as akin to regularization

$$\hat{\theta}^{MAP} = \arg \max_{\theta} \log p(\theta | x)$$

- Saw importance of regularization in logistic regression (ML estimate can overfit data and lead to poor generalization)

©Emily Fox 2013

36

EM Algorithm – MAP Case

- Re-derive EM algorithm for $p(\theta | x)$
- Add $\log p(\theta)$ to $U(\theta, \hat{\theta}^{(t)})$
 - What must be computed in E-Step remains unchanged because this term does not depend on y .
 - M-Step becomes:

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

©Emily Fox 2013

37

MAP EM Example – MoG

- For mixture of Gaussians, conjugate priors are:
 $\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad \{\mu_k, \Sigma_k\} \sim \text{NIW}(m_0, \kappa_0, \nu_0, S_0)$
- Results in following M-Step:

$$\hat{\mu}_k = \frac{r_k \bar{x}_k + \kappa_0 m_0}{r_k + \kappa_0} \quad \hat{\pi}_k = \frac{r_k + \alpha_k - 1}{N + \sum_k \alpha_k - K}$$

$$\hat{\Sigma}_k = \frac{S_0 + r_k S_k + \frac{\kappa_0 r_k}{\kappa_0 + r_k} (\bar{x}_k - m_0)(\bar{x}_k - m_0)'}{\nu_0 + r_k + d + 2}$$

©Emily Fox 2013

38

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

Course Announcements

- Homework 2 will be posted on Thursday
 - Due 2 weeks later (Feb 14)
- Project proposals:
 - Initial ideas now posted
 - Deadline extended to Tues, Feb 5
 - 1 page, 1-2 people
- Recitation on Thursday (Linda)
- Office hours as normal