

Case Study 2: Document Retrieval

Review: Mixtures of Gaussians

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

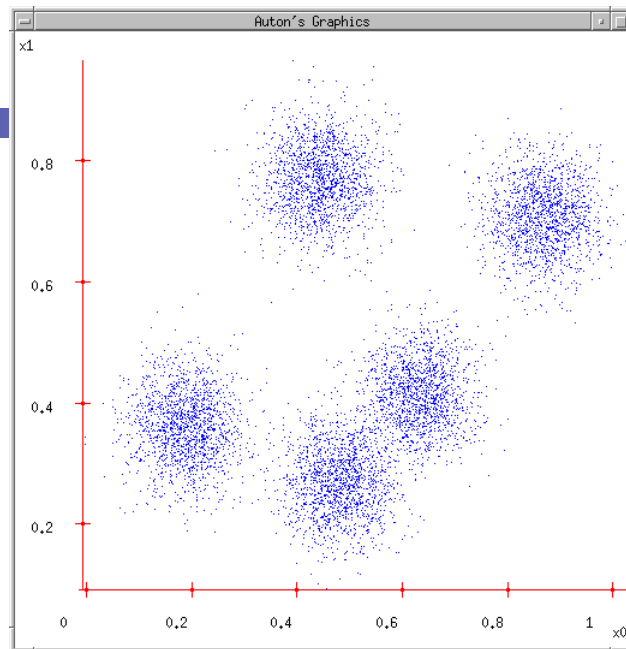
Emily Fox

January 28th, 2014

©Emily Fox 2014

1

Some Data



©Emily Fox 2014

2

Gaussian Mixture Model

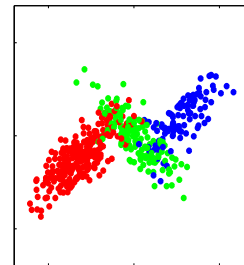
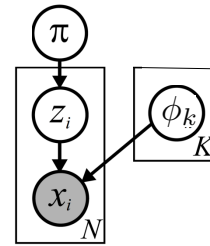
- Most commonly used mixture model
- Observations:

- Parameters:

- Cluster indicator:

- Per-cluster likelihood:

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



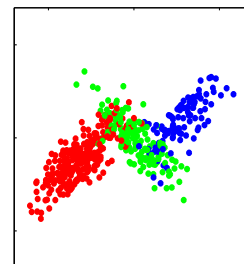
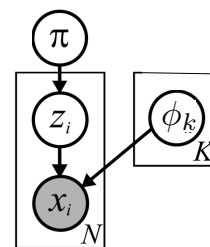
©Emily Fox 2014

3

Generative Model

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

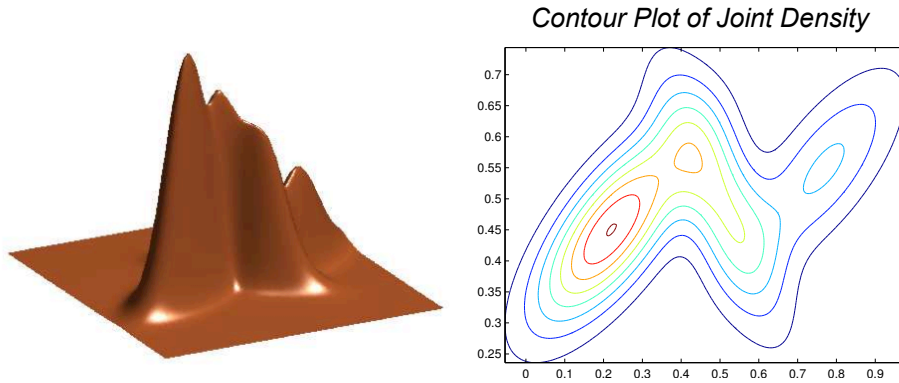
- For each observation i ,
 - Sample a cluster assignment
 - Sample the observation from the selected Gaussian



©Emily Fox 2014

4

Also Useful for Density Estimation

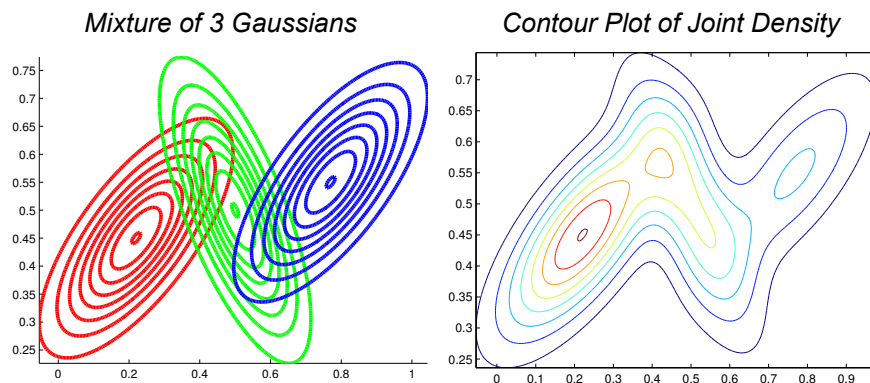


©Emily Fox 2014

5

Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians



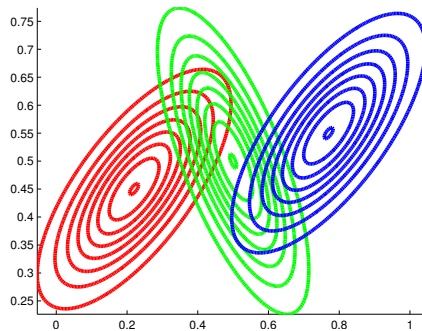
©Emily Fox 2014

6

Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



$$p(x^i | \pi, \mu, \Sigma) =$$

©Emily Fox 2014

7

Summary of GMM Components

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

Gaussian mixture marginal and conditional likelihood :

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

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

©Emily Fox 2014

8

Document Representation

- Bag of words model



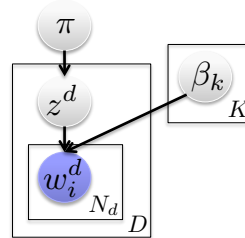
document d

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 2014

13

Form of Likelihood

- Conditioned on topic...

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

- Marginalizing latent topic assignment:

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

©Emily Fox 2014

14

Case Study 2: Document Retrieval

Review: EM Algorithm

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

Emily Fox
January 28th, 2014

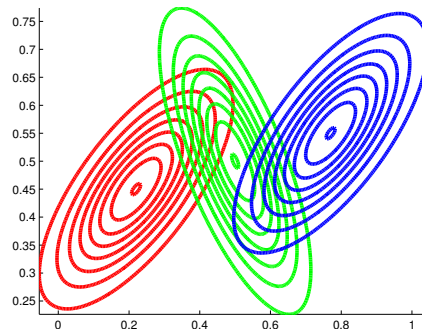
©Emily Fox 2014

15

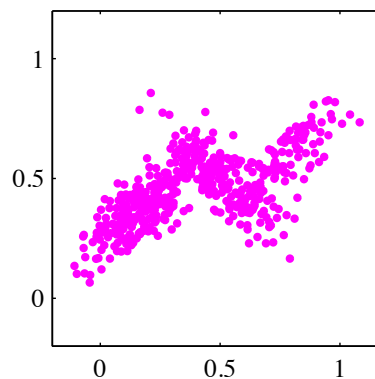
Learning Model Parameters

- Want to learn model parameters

Mixture of 3 Gaussians



Our actual observations



C. Bishop, *Pattern Recognition & Machine Learning*

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

$$L_x(\theta) =$$

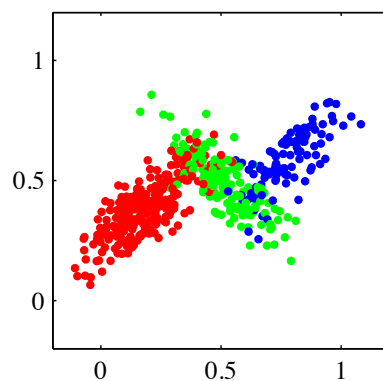
- Neither convex nor concave and local optima

©Emily Fox 2014

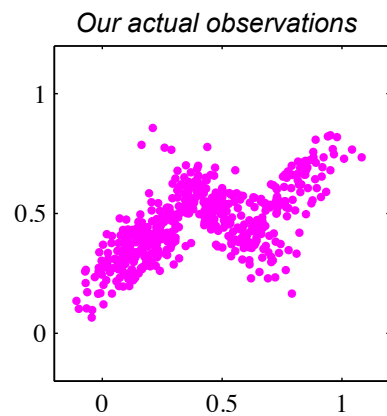
17

Complete Data

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



Complete data labeled
by true cluster assignments



Our actual observations

C. Bishop, *Pattern Recognition & Machine Learning*

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 $k!$

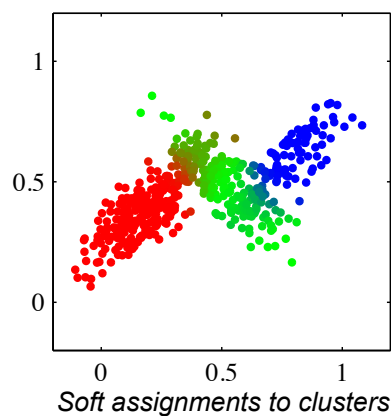
- Example: mixture of Gaussians (MoG) $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

©Emily Fox 2014

19

Cluster Responsibilities

- We must infer the cluster assignments from the observations



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

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

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

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

1. Infer “responsibilities”

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

2. Optimize parameters

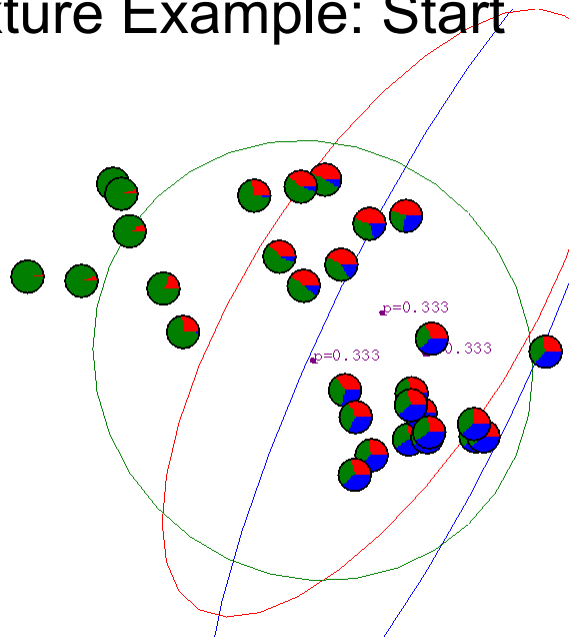
max w.r.t. π_k :

max w.r.t. ϕ_k :

©Emily Fox 2014

21

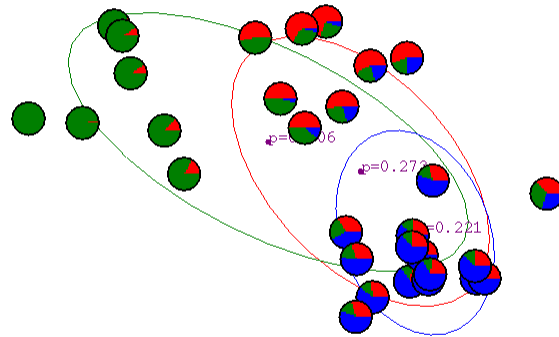
Gaussian Mixture Example: Start



©Emily Fox 2014

22

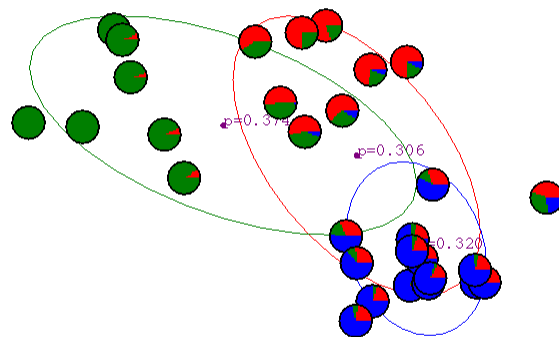
After first iteration



©Emily Fox 2014

23

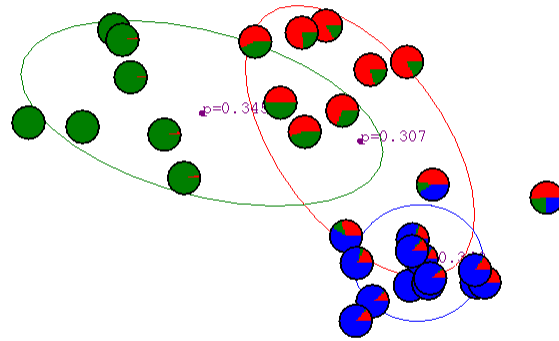
After 2nd iteration



©Emily Fox 2014

24

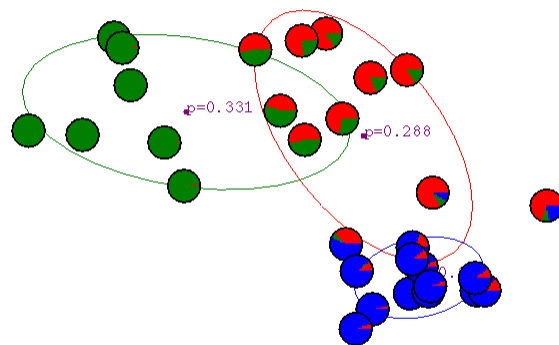
After 3rd iteration



©Emily Fox 2014

25

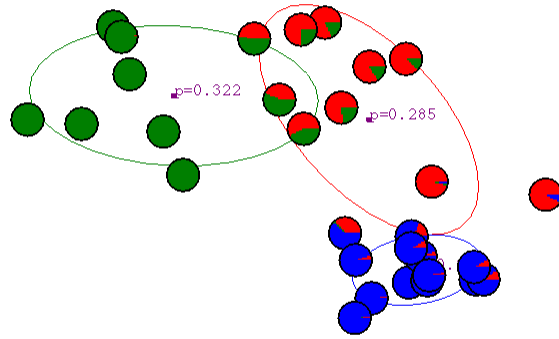
After 4th iteration



©Emily Fox 2014

26

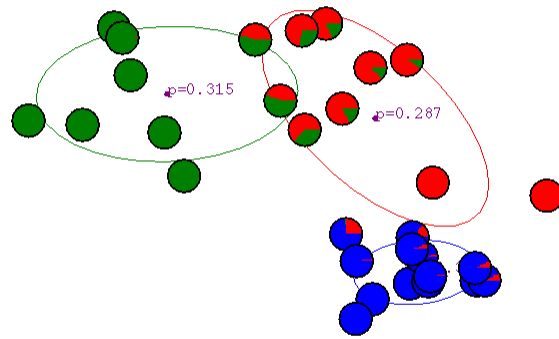
After 5th iteration



©Emily Fox 2014

27

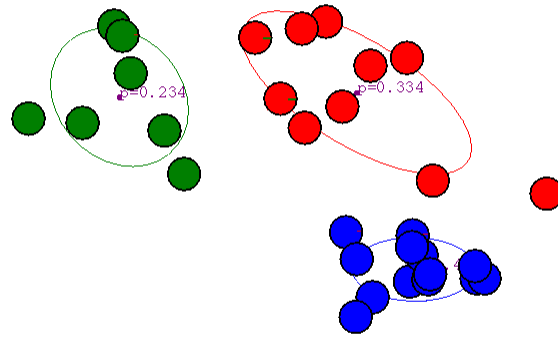
After 6th iteration



©Emily Fox 2014

28

After 20th iteration



©Emily Fox 2014

29

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 2014

30

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.

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

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

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

33

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 2014

34

Case Study 2: Document Retrieval

Review: Connection to k-means

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

Emily Fox

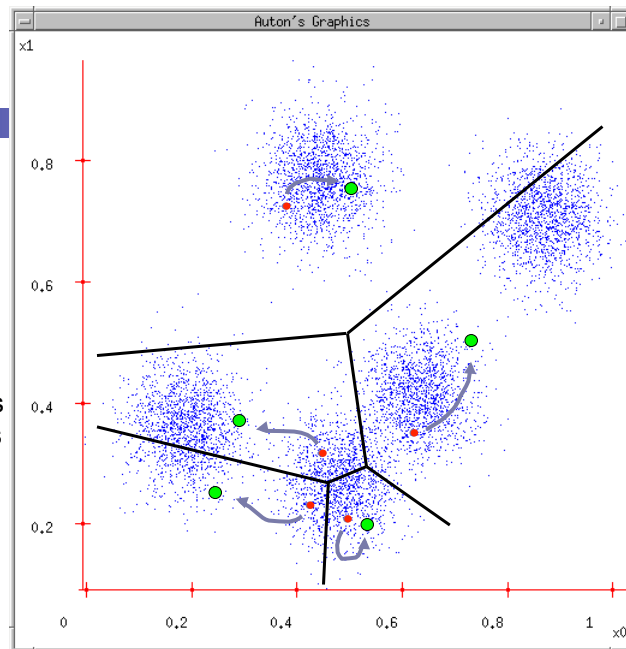
January 28th, 2014

©Emily Fox 2014

35

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



©Emily Fox 2014

36

K-means

- Randomly initialize k centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $z^j \leftarrow \arg \min_i \|\mu_i - \mathbf{x}^j\|_2^2$

- **Recenter:** μ_i becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: z^j = i} \|\mu - \mathbf{x}^j\|_2^2$

- Equivalent to $\mu_i \leftarrow$ average of its points!

©Emily Fox 2014

37

Special case: spherical Gaussians and 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:

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

- Then, compare EM objective with k-means:

©Emily Fox 2014

38