

So far, supervised learning

$h: X \rightarrow \mathbb{R}$  "regression"

$h: X \rightarrow \{0, 1, \dots, K\}$   
"classification"

Unsupervised

# Clustering

## K-means

Machine Learning – CSE546

Emily Fox

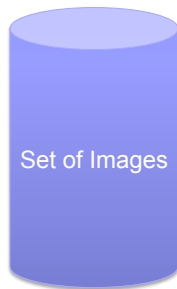
University of Washington

November 4, 2013

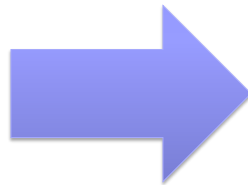
©Carlos Guestrin 2005-2013

1

## Clustering images

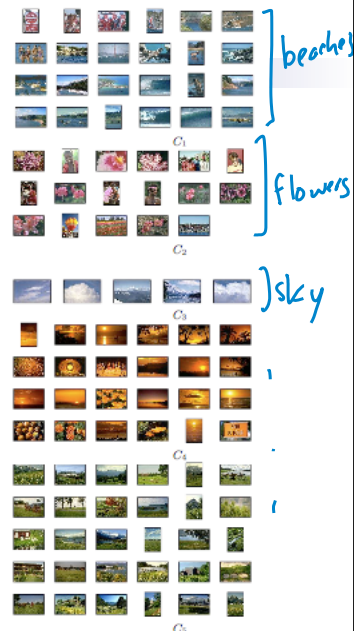


Set of Images



organize into  
coherent "themes"

key: no labels given

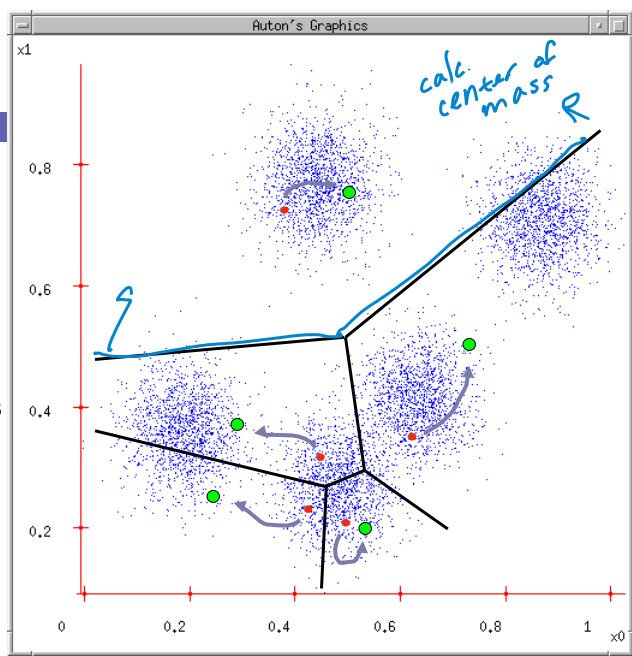


©Carlos Guestrin 2005-2013

[Goldberger et al.] 2

# 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



©Carlos Guestrin 2005-2013

3

# K-means

- Randomly initialize  $k$  centers (or "smartly")
  - $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$  *iteration*
- **Classify:** Assign each point  $j \in \{1, \dots, N\}$  to nearest center:
  - $C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$  *fix  $\mu$ , opt.  $C$*
  - $C(j)=k \Rightarrow j$ th obs. is assoc. w/ cluster  $k$*
- **Recenter:**  $\mu_i$  becomes centroid of its point: *fix  $C$ , opt.  $\mu$* 
  - $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j:C(j)=i} \|\mu - x_j\|^2 \leftarrow \mu_i = \frac{\sum_{j:C(j)=i} x_j}{|\{j : C(j)=i\}|}$
  - Equivalent to  $\mu_i \leftarrow$  average of its points!  *$\{j : C(j)=i\}$*

©Carlos Guestrin 2005-2013

4

*model that can be used for clustering, density est.:*

# Mixtures of Gaussians

Machine Learning – CSE546  
 Emily Fox  
 University of Washington

November 4, 2013  
©Carlos Guestrin 2005-2013

5

## (One) bad case for k-means

- Clusters may overlap
- Some clusters may be "wider" than others

*shape*

*centers*

*params defining clusters*

*so centers alone don't tell the whole story*

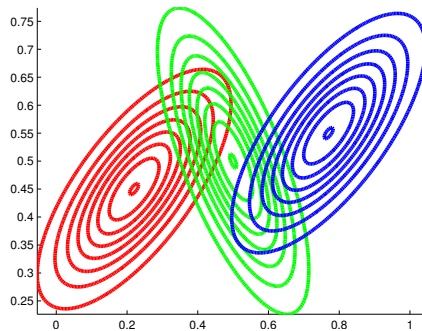
©Carlos Guestrin 2005-2013

6

# Density as Mixture of Gaussians

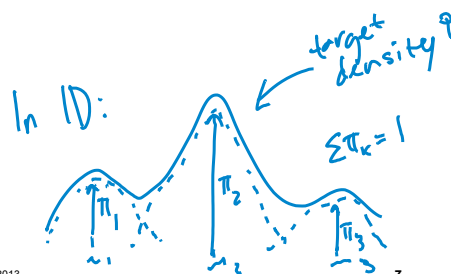
- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



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

*Handwritten notes:*  $\pi = [\pi_1, \dots, \pi_K]$ ,  $\{\mu_k, \Sigma_k\}$

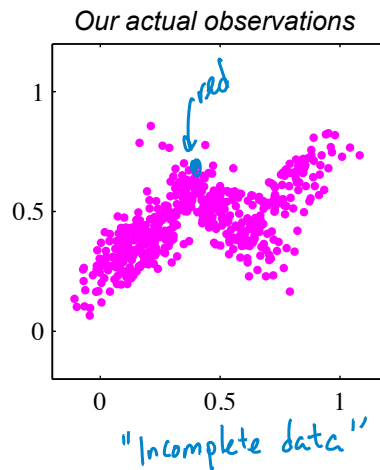
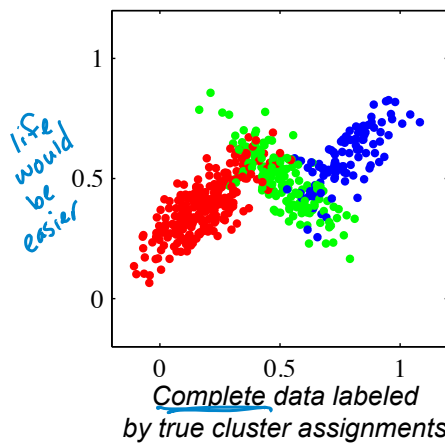


©Emily Fox 2013

7

# Clustering our Observations

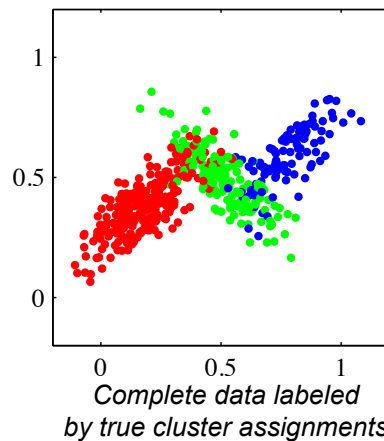
- Imagine we have an assignment of each  $x^i$  to a Gaussian



C. Bishop, Pattern Recognition & Machine Learning

# Clustering our Observations

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



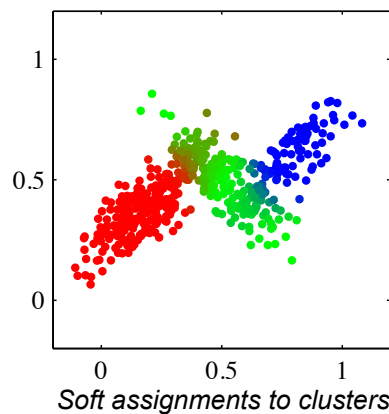
- Introduce latent cluster indicator variable  $z^i$  ( $i$ )  $\rightarrow$   $z^i$   
 $z^i \in \{1, \dots, K\} \equiv$   
 $\Pr(z^i = k) = \pi_k$
- Then we have  
 $p(x^i | z^i = k, \mu, \Sigma) = N(x^i | \mu_k, \Sigma_k)$

param est. is easy if we have  $\{z^i\}$   
 $\Rightarrow$  decouples into  $K$  Gauss. est.

C. Bishop, Pattern Recognition & Machine Learning

# Clustering our Observations

- We must infer the cluster assignments from the observations "responsibilities"



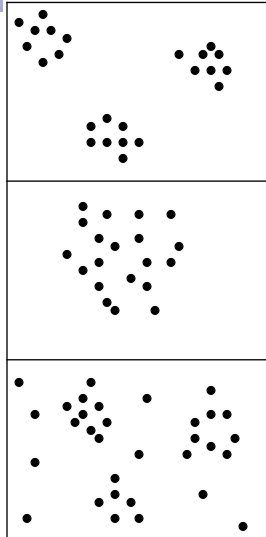
- Posterior probabilities of assignments to each cluster \*given\* model parameters:  
 $r_{ik} = p(z^i = k | x^i, \pi, \mu, \Sigma) =$

$$= \frac{\pi_k N(x^i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x^i | \mu_j, \Sigma_j)}$$

motivates an iterative alg.

C. Bishop, Pattern Recognition & Machine Learning

# Unsupervised Learning: not as hard as it looks



Sometimes easy

Sometimes impossible

and sometimes in between

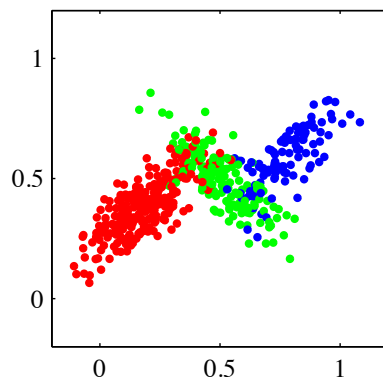
©Carlos Guestrin 2005-2013

11

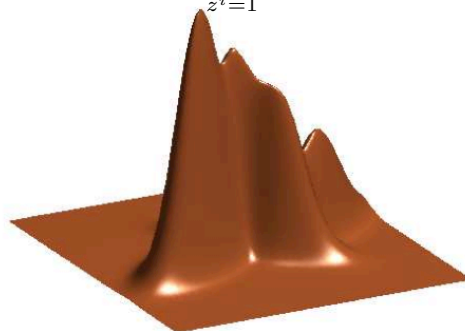
# Summary of GMM Concept

- Estimate a density based on  $x^1, \dots, x^N$

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



Complete data labeled  
by true cluster assignments



Surface Plot of Joint Density,  
Marginalizing Cluster Assignments

©Emily Fox 2013

12

## 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 :**

$$p(x^i | \pi, \mu, \Sigma) = \sum_{z^i=1}^K \pi_{z^i} p(x^i | z^i, \mu, \Sigma)$$

$$p(x^i | z^i, \mu, \Sigma) = \mathcal{N}(x^i | \mu_{z^i}, \Sigma_{z^i})$$

©Emily Fox 2013

13

## Expectation Maximization

Machine Learning – CSE546

Emily Fox

University of Washington

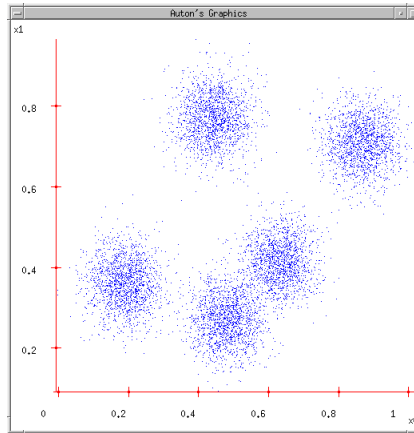
November 6, 2013

©Carlos Guestrin 2005-2013

14

Next... back to Density Estimation

What if we want to do density estimation with multimodal or clumpy data?



©Carlos Guestrin 2005-2013

15

But we don't see class labels!!!

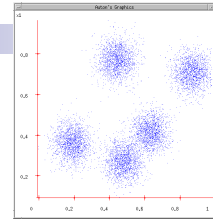
■ MLE:

□  $\operatorname{argmax} \prod_i P(z^i, x^i)$

■ But we don't know  $z^i$

■ Maximize marginal likelihood:

□  $\operatorname{argmax} \prod_i P(x^i) = \operatorname{argmax} \prod_i \sum_{k=1}^K P(z^i=k, x^i)$



©Carlos Guestrin 2005-2013

16



## Special case: spherical Gaussians and hard assignments

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

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

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

- If each  $\mathbf{x}^i$  belongs to one class  $C(i)$  (hard assignment), marginal likelihood:

$$\prod_{i=1}^N \sum_{k=1}^K P(\mathbf{x}^i, z^i = k) \propto \prod_{i=1}^N \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_{C(i)}\|^2\right]$$

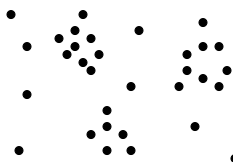
- Same as K-means!!!

©Carlos Guestrin 2005-2013

17

## EM: “Reducing” Unsupervised Learning to Supervised Learning

- If we knew assignment of points to classes → Supervised Learning!



- Expectation-Maximization (EM)

- Guess assignment of points to classes

- In standard (“soft”) EM: each point associated with prob. of being in each class

- Recompute model parameters

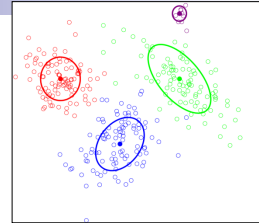
- Iterate

©Carlos Guestrin 2005-2013

18

# Generic Mixture Models

MoG Example:



- Observations:
- Parameters:
- Likelihood:
- Ex.  $z^i$  = country of origin,  $x^i$  = height of  $i^{\text{th}}$  person
  - $k^{\text{th}}$  mixture component = distribution of heights in country  $k$

©Emily Fox 2013

19

# 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

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

- Neither convex nor concave and local optima

©Emily Fox 2013

20

## 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)  $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$

## 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”

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

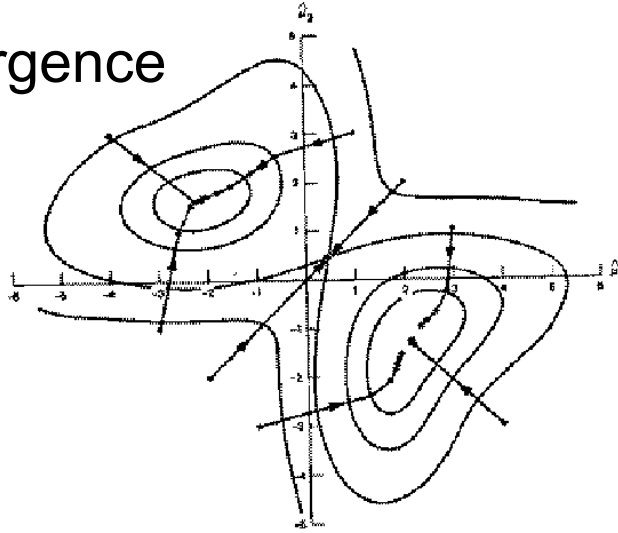
2. Optimize parameters

max w.r.t.  $\pi_k$  :

max w.r.t.  $\mu_k, \Sigma_k$  :

# E.M. Convergence

- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func.  $\rightarrow$  convergence to a local optimum guaranteed

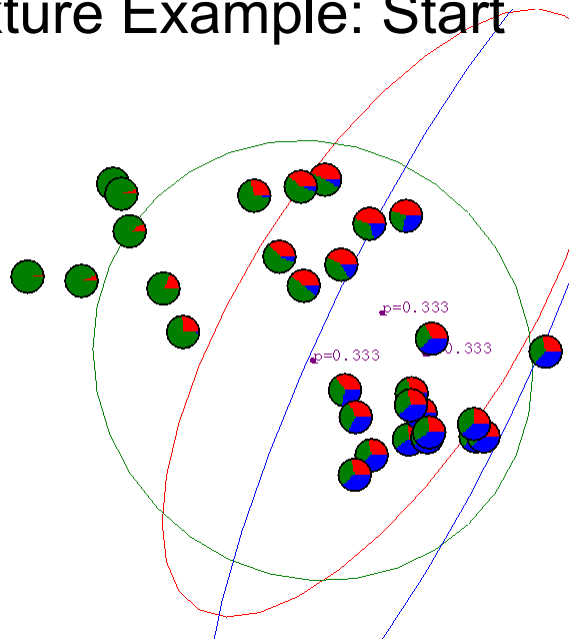


- This algorithm is REALLY USED. And in high dimensional state spaces, too. E.G. Vector Quantization for Speech Data

©Carlos Guestrin 2005-2013

23

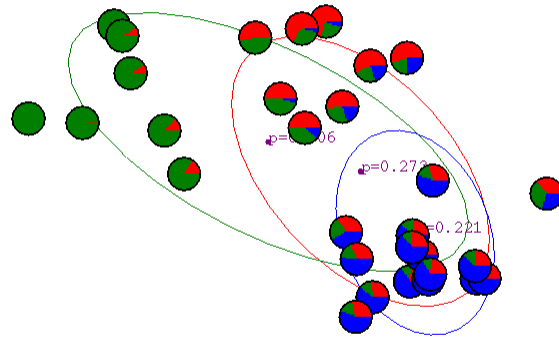
# Gaussian Mixture Example: Start



©Emily Fox 2013

24

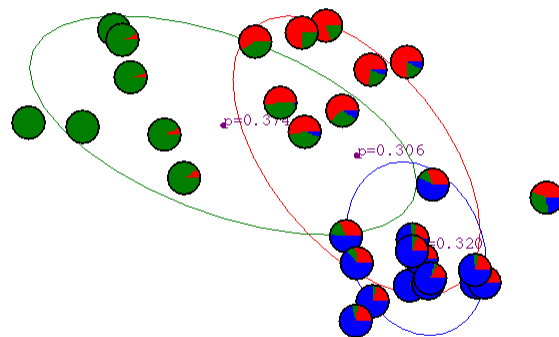
# After first iteration



©Emily Fox 2013

25

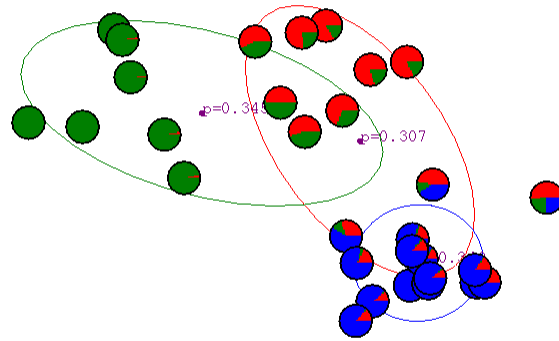
# After 2nd iteration



©Emily Fox 2013

26

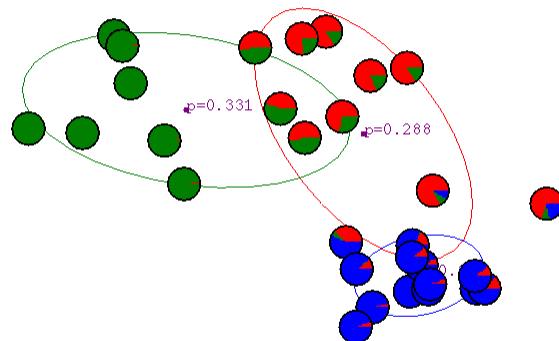
## After 3rd iteration



©Emily Fox 2013

27

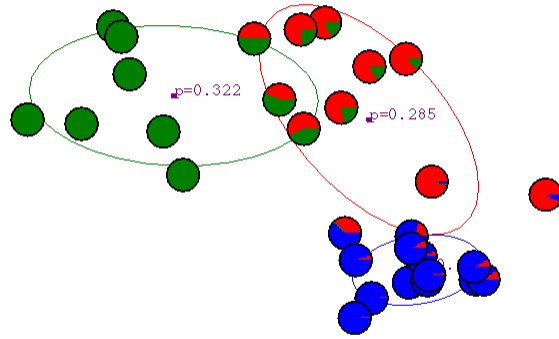
## After 4th iteration



©Emily Fox 2013

28

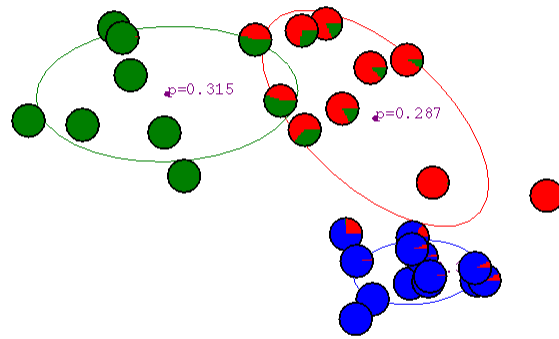
## After 5th iteration



©Emily Fox 2013

29

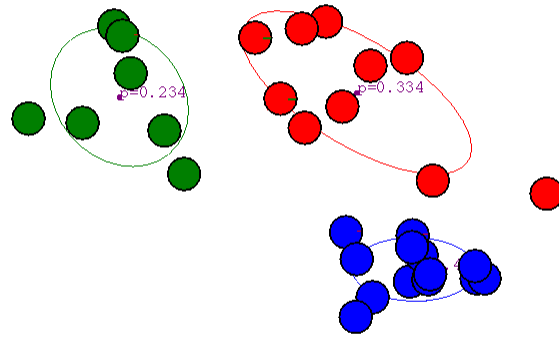
## After 6th iteration



©Emily Fox 2013

30

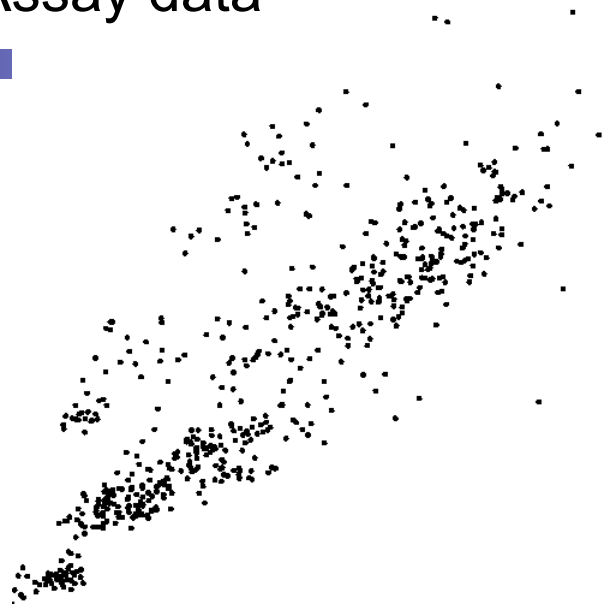
## After 20th iteration



©Emily Fox 2013

31

## Some Bio Assay data

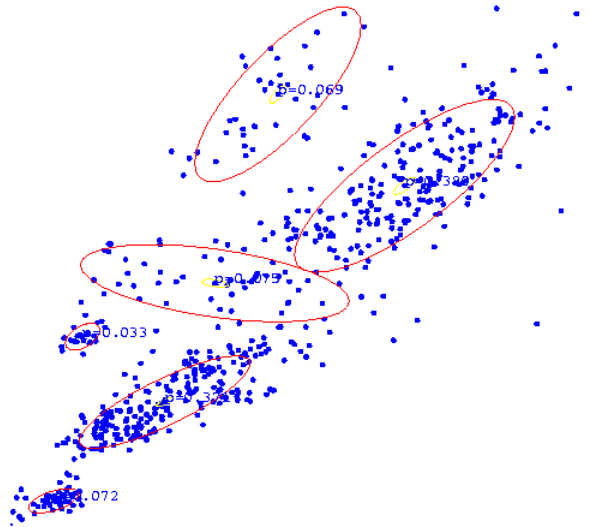


©Emily Fox 2013

32



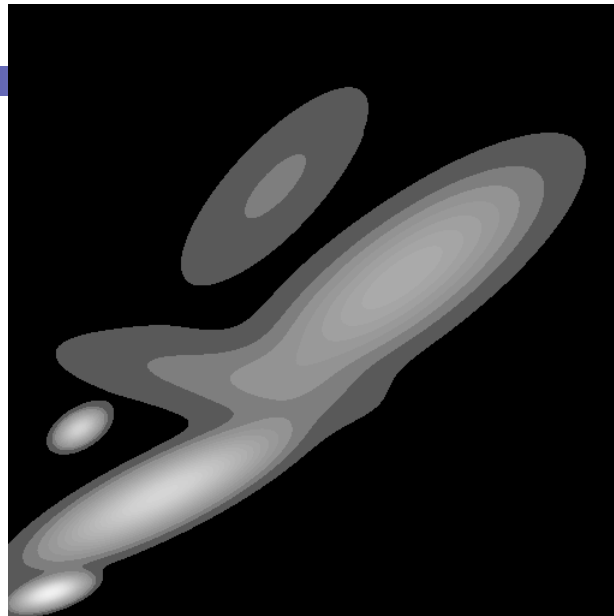
# GMM clustering of the assay data



©Emily Fox 2013

33

## Resulting Density Estimator



©Emily Fox 2013

34

## 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  
 $\theta$  parameters

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

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

- Special case:

$$x = g(y)$$

©Emily Fox 2013

35

## 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

36

## Expectation Maximization (EM) – Derivation

- Step 3

- Consider log likelihood of data at any  $\theta$  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  $\theta$  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:

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

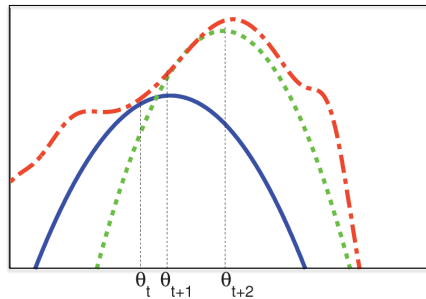


Figure from  
KM textbook

©Emily Fox 2013

41

## Comments on EM

- Since Gibbs inequality is satisfied with equality only if  $p=q$ , any step that changes  $\theta$  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

42

## 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

43

## What you should know

- K-means for clustering:
  - algorithm
  - converges because it’s coordinate ascent
- EM for mixture of Gaussians:
  - How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent

©Carlos Guestrin 2005-2013

44