So far, supervised learning    h: X → R  "regression"
h: X → {0,1...,K}
"classification"

Unsupervised

# Clustering
# K-means

Machine Learning – CSE546

Emily Fox

University of Washington

November 4, 2013

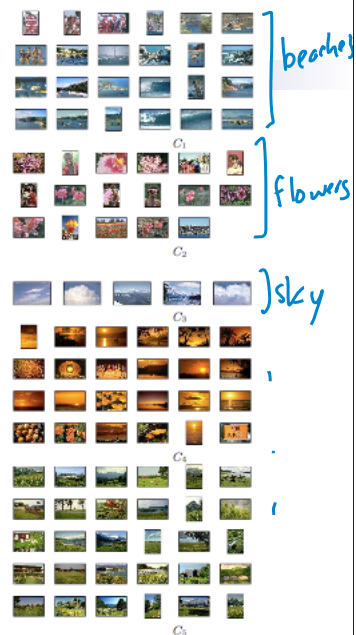©Carlos Guestrin 2005-2013

1

---

# Clustering images

key: no labels given

beaches

flowers

sky

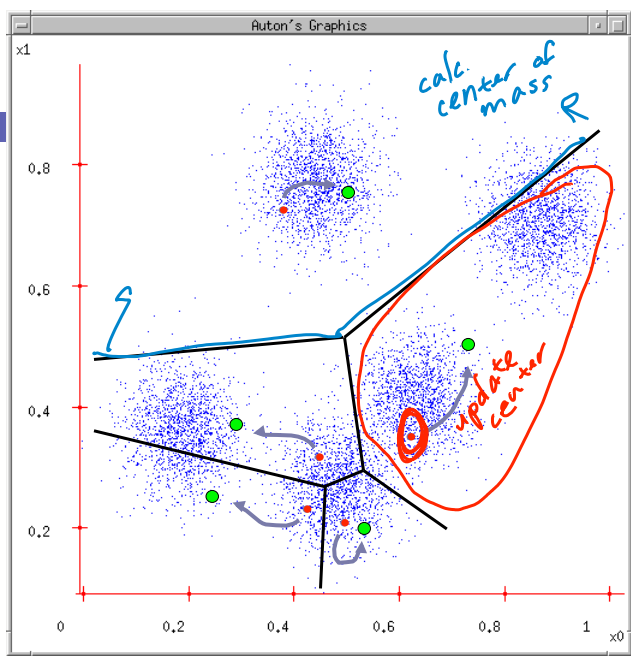Set of Images

organize into coherent "themes"

©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

*calc. center of mass*

*update center*

3

---

# K-means

*Coord. desc. alg. → converges to local mode*

- Randomly initialize *k* centers    *(or "smartly")*
  - $\mu^{(0)} = \mu_1^{(0)},\ldots, \mu_k^{(0)}$   *iteration*

*converged when nothing moves (no point changes its cluster)*

- **Classify**: Assign each point $j \in \{1,\ldots N\}$ to nearest center:
  - $C^{(t)}(j) \leftarrow \arg\min_i ||\mu_i - x_j||^2$

*ith cluster center*  *jth data pt*  *fix μ, opt. C*

*C(j)=k ⇒ jth obs. is assoc. w/ cluster k*

*iterate.*

- **Recenter**: $\mu_i$ becomes centroid of its point:   *fix C, opt. μ*
  - $\mu_i^{(t+1)} \leftarrow \arg\min_\mu \sum_{j:C(j)=i} ||\mu - x_j||^2$   $\leftarrow \mu_i = \frac{\sum_{j:C(j)=i} x_j}{|\{j:C(j)=i\}|}$
  - Equivalent to $\mu_i \leftarrow$ average of its points!

4

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

# Mixtures of Gaussians

Machine Learning – CSE546

Emily Fox

University of Washington

November 4, 2013

5

---

# (One) bad case for k-means

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

*shape*

*params defining clusters*
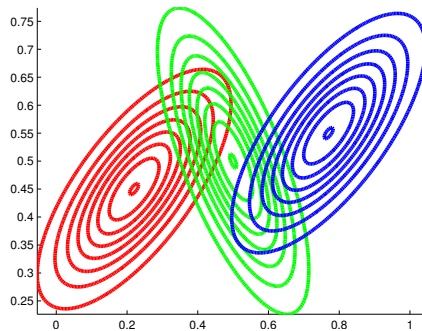
*centers*

*so centers alone don't tell the whole story*
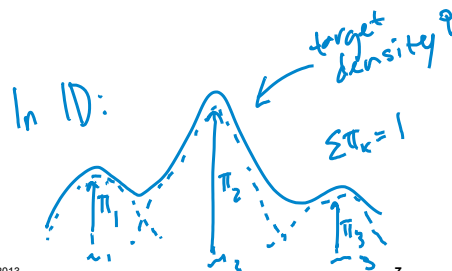
6

3

# Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

*Mixture of 3 Gaussians*



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

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

$[\pi_1, \ldots, \pi_K]$   $\{\mu_k, \Sigma_k\}$

In 1D:   target density $q$   $\sum \pi_k = 1$

$\pi_1$   $\pi_2$   $\pi_3$

©Emily Fox 2013   7

# Clustering our Observations

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

*Our actual observations*



life would be easier

red

*Complete data labeled by true cluster assignments*

"Incomplete data"

*C. Bishop, Pattern Recognition & Machine Learning*

©Emily Fox 2013

# Clustering our Observations

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



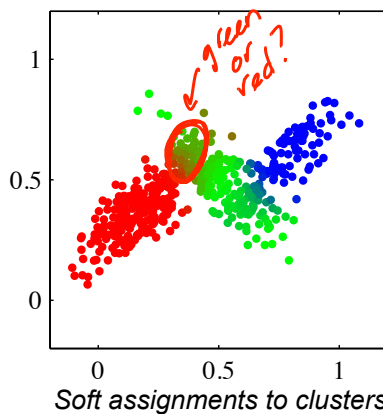*Complete data labeled by true cluster assignments*

- Introduce latent cluster indicator variable $z^i$

  $c(i) \rightarrow z^i$

  $z^i \in \{1, ..., k\}$

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

- Then we have

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

  param est. is easy if we have $\{z_i\}$

  $\Rightarrow$ decouples into $k$ Gauss. est.

*C. Bishop, Pattern Recognition & Machine Learning*

©Emily Fox 2013


# Clustering our Observations

- We must infer the cluster assignments from the observations

  "responsibilities"



*Soft assignments to clusters*

  green or red?

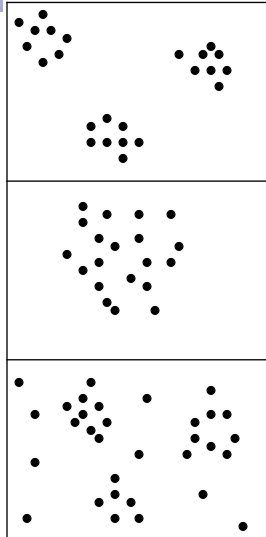- Posterior probabilities of assignments to each cluster *given* model parameters:

  resp.

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

  obs $i$   cluster $k$   $= \dfrac{\pi_k N(x^i | \mu_k, \Sigma_k)}{\sum\limits_{j=1}^{K} \pi_j N(x^i | \mu_j, \Sigma_j)}$

  motivates an iterative alg.

*C. Bishop, Pattern Recognition & Machine Learning*

©Emily Fox 2013

# Unsupervised Learning:
## not as hard as it looks



Sometimes easy

Sometimes impossible

and sometimes in between
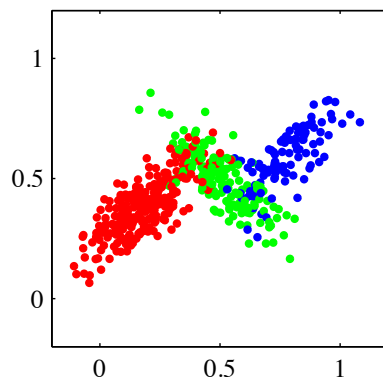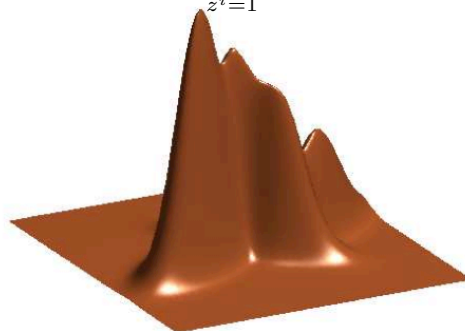
---

# Summary of GMM Concept

■ Estimate a density based on $x^1, \ldots, x^N$

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



*Complete data labeled
by true cluster assignments*

*Surface Plot of Joint Density,
Marginalizing Cluster Assignments*

# Summary of GMM Components

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

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

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

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

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

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

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

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

13

---

*iterative alg. for MLE*

# Expectation Maximization

Machine Learning – CSE546

Emily Fox

University of Washington

November 6, 2013

14
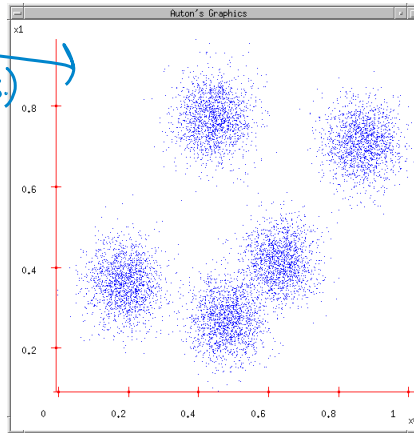
# Next… back to Density Estimation

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



*Goal is to fit a MoG (mix. of Gauss.) to this data*

*Learn: $\{\pi_k, \mu_k, \Sigma_k\}$*

15

---

# But we don't see class labels!!!

*In classification $x^i = \{GPA=3.9, ML\ Grade=4.0,...\}$*

*$z^i = \{Role = VP\}$*



- **MLE:**
  - argmax $\prod_i P(z^i, x^i)$
    - $\pi, \mu, \Sigma$ — *class labels* — *features*

*"nuissance variable"*

- But we don't know $z^i$ ←

- Maximize <u>marginal likelihood</u>:
  - argmax $\prod_i P(x^i)$ = argmax $\prod_i \sum_{k=1}^{K} P(z^i=k, x^i)$
    - $\pi, \mu, \Sigma$ — *only obs. this*
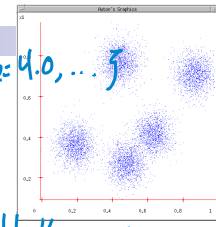
*$N(x^i; \mu_k, \Sigma_k)$*
*$P(x^i | z^i=k)$*
*$P(z^i=k)$*
*$\pi_k$*

*Sum/avg. out unobserved variables*

*Sum role = {VP, Engineer, Barista}*
*weigh by prob.*
*$P(z=VP), P(z=Eng.), P(z=barista)$*

16

8

# Special case: spherical Gaussians and hard assignments

*Handwritten annotations:* $P(x^i | z^i = k)$  $P(z^i = k)$

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

*Handwritten:* $x^i \in \mathbb{R}^d$

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

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

*Handwritten:* $\Sigma_k = \begin{bmatrix} \sigma^2 & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix} = \sigma^2 I$

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

*Handwritten:* $\Rightarrow P(z^i = k) = \begin{cases} 1 & C(i) = k \\ 0 & \text{otherwise} \end{cases}$

*Handwritten:* want to max this:

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

*Handwritten:* marginal like.

- Same as K-means!!!

*Handwritten:*
$$\max_{C, \mu, \sigma} \prod_{i=1}^{} e^{-\frac{1}{2\sigma^2} \| x^i - \mu_{C(i)} \|^2} = \max \ln \prod = \max \sum -\frac{1}{2\sigma^2} \| x^i - \mu_{C(i)} \|^2$$
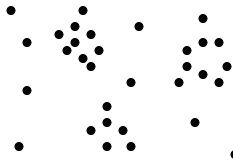$$= \min_{C, \mu, \sigma} \sum_{i=1}^{} \| x^i - \mu_{C(i)} \|^2 \leftarrow \text{exactly k-means obj.}$$

©Carlos Guestrin 2005-2013    17

---

# EM: "Reducing" Unsupervised Learning to Supervised Learning

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

  *Handwritten:* easy

- Expectation-Maximization (EM)
  - Guess assignment of points to classes
    - In standard ("soft") EM: each point associated with prob. of being in each class
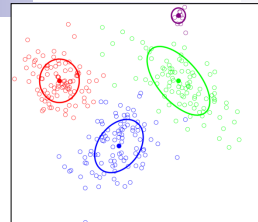  - Recompute model parameters
  - Iterate

  *Handwritten:* iterate until convergence like in k-means

©Carlos Guestrin 2005-2013    18

9

# Generic Mixture Models

*MoG Example:*

- Observations: $x^1, \ldots, x^N$ with $x^i \in \mathbb{R}^d$

- Parameters:
  $$\pi = [\pi_1, \ldots, \pi_K] \quad \text{mix. weights}$$
  $$\phi = \{\phi_1, \ldots, \phi_k\] \quad \text{like. params,}$$
  (defines cluster 1) (e.g. $\phi_k = \{\mu_k, \Sigma_k\}$ for MoG)
  $$\theta = \{\pi, \phi\}$$

- Likelihood:
  $$p(x^i | \theta) = \sum_{k=1}^{K} \pi_k \, p(x^i | \phi_k) \qquad \text{e.g. } N(x^i, \mu_k, \Sigma_k)$$

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

19

---

# ML Estimate of Mixture Model Params

Marginal

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

  $$P(X|\theta) = \prod_i p(x^i | \theta) = \prod_i \sum_{z^i} p(x^i, z^i | \theta)$$

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

- Neither convex nor concave and local optima

20

10

## If "complete" data were observed...

$$\sum_i \log \pi_{z_i} = \sum_{j=1}^{K} \sum_{i: z_i = j} \log \pi_j = \sum_{j=1} N_j \log \pi_j \qquad N_j = |\{i: z_i = j\}|$$

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

$$L_{x,z}(\theta) = \sum_i \log p(x^i, z^i \mid \theta) = \sum_i \log p(x^i \mid z^i, \theta) + \log p(z^i \mid \theta) \quad \overbrace{\pi_{z^i}}$$

$$= \sum_{j=1}^{K} \sum_{i: z_i = j} \log p(x^i \mid z^i = j, \phi_j) + \sum_{j=1}^{K-1} N_j \log \pi_j + N_K \log\left(1 - \sum_{j=1}^{K-1} \pi_j\right)$$

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

$$\hat{\phi}_k = \arg\max_{\phi_k} \sum_{i: z_i = k} \log p(x^i \mid z^i = k, \phi_k) \qquad \hat{\pi}_k = \frac{N_k}{N} \quad k = 1, \ldots, K-1$$

$$\sum \pi_j = 1$$

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

$$\begin{cases} \hat{\mu}_k = \dfrac{1}{N_k} \sum_{i: z^i = k} x_i \\[2mm] \hat{\Sigma}_k = \dfrac{1}{N_k} \sum_{i: z^i = k} x_i x_i^T - \hat{\mu}_k \hat{\mu}_k^T \end{cases} \qquad \hat{\pi}_k = \frac{N_k}{N}$$

©Emily Fox 2013    21

---

## 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

  $$z^i \rightleftharpoons \hat{\theta} \quad \text{est.} \qquad \hat{\theta} = \{\phi_k, \pi_k\}$$

- Example: MoG (derivation soon... + HW)
  1. Infer "responsibilities"  $\quad$ prev. iter.

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

  2. Optimize parameters

  $$\max \text{ w.r.t. } \pi_k : \quad \hat{\pi}_k^{(t)} = \frac{1}{N} \sum_i r_{ik} = \frac{r_k}{N} \leftarrow \text{soft counts!}$$
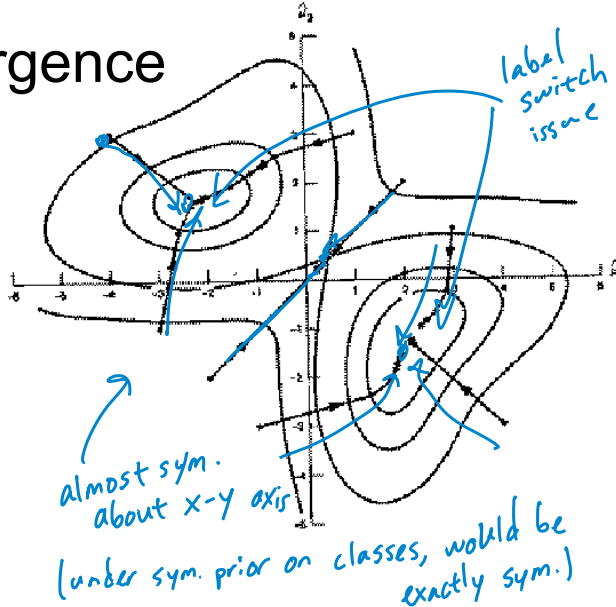
  $$\max \text{ w.r.t. } \mu_k, \Sigma_k :$$

  $$\hat{\mu}_k^{(t)} = \frac{\sum_i r_{ik} x_i}{N} \leftarrow \text{weighted mean} \qquad \hat{\Sigma}_k^{(t)} = \frac{1}{r_k} \sum r_{ik} x_i x_i^T - \hat{\mu}_k^{(t)} \hat{\mu}_k^{(t)T}$$

©Emily Fox 2013    22

11

# E.M. Convergence

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

*[handwritten annotations: "label switch issue", "almost sym. about x-y axis", "(under sym. prior on classes, would be exactly sym.)"]*

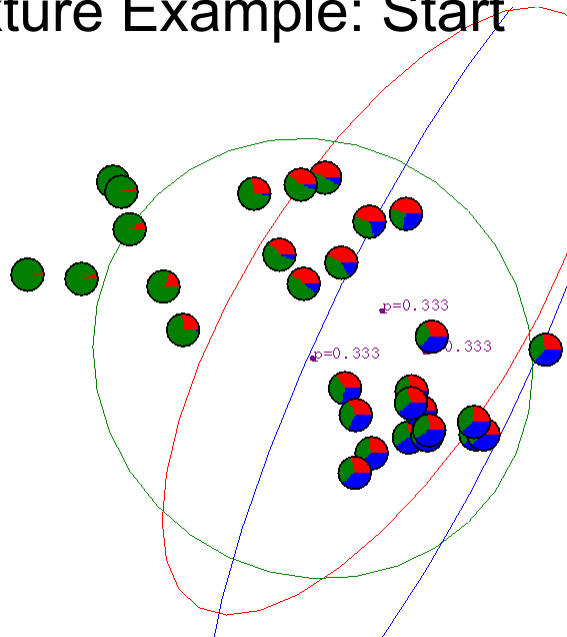■ This algorithm is REALLY USED. And in high dimensional state spaces, too. E.G. Vector Quantization for Speech Data

23

# Gaussian Mixture Example: Start

*[handwritten annotations: "start with initial est. of $\pi^{(0)}, \theta^{(0)}$", "→ lead to initial "responsibilities""]*

p=0.333
p=0.333    0.333

24

12
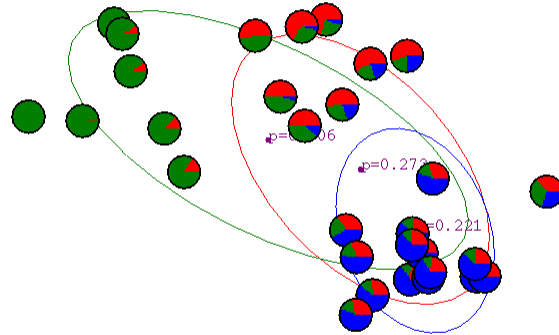
# After first iteration

max. like
given soft
assignments

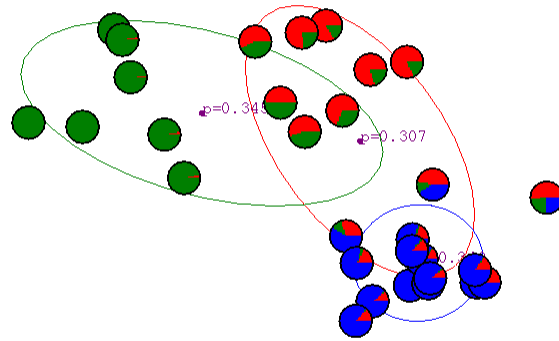→ use new
$\pi^{(1)}, \phi^{(1)}$
to compute
new $r_{ik}$

25

# After 2nd iteration

26

13

# After 3rd iteration
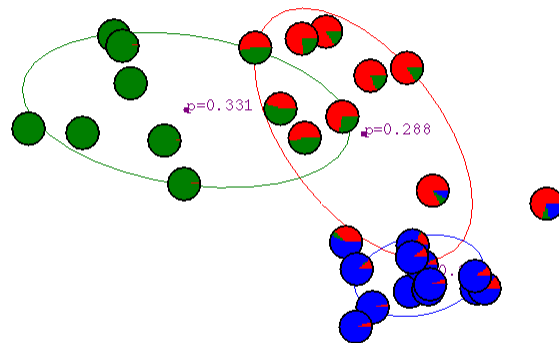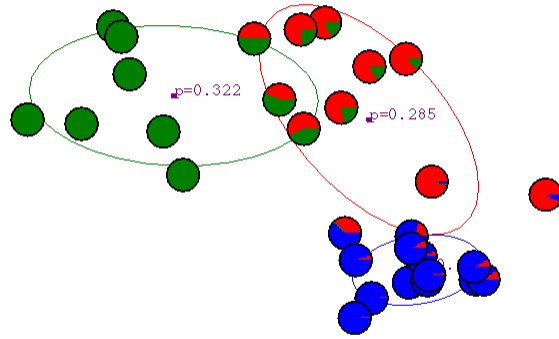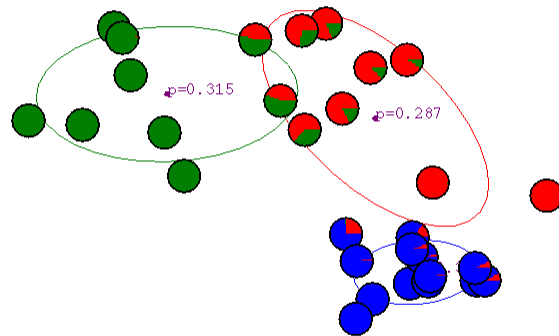
27

# After 4th iteration

28

14

# After 5th iteration



©Emily Fox 2013

29

# After 6th iteration



©Emily Fox 2013

30

15

# After 20th iteration



*looks pretty good*

31

# Some Bio Assay data

32

16

# GMM clustering of the assay data



p=0.069
p=0.075
0.033
.072

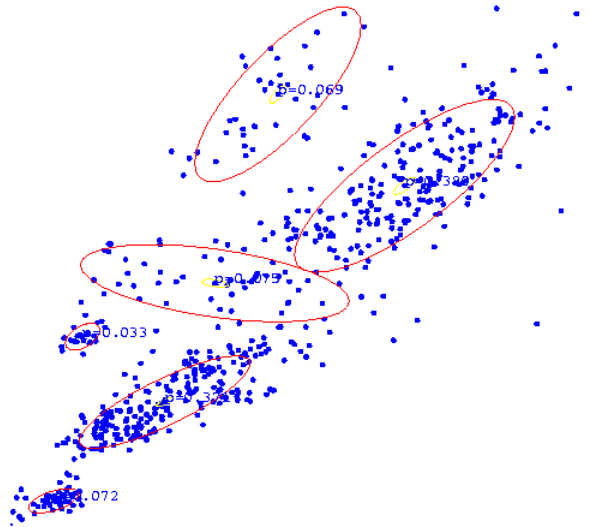©Emily Fox 2013                                                                33
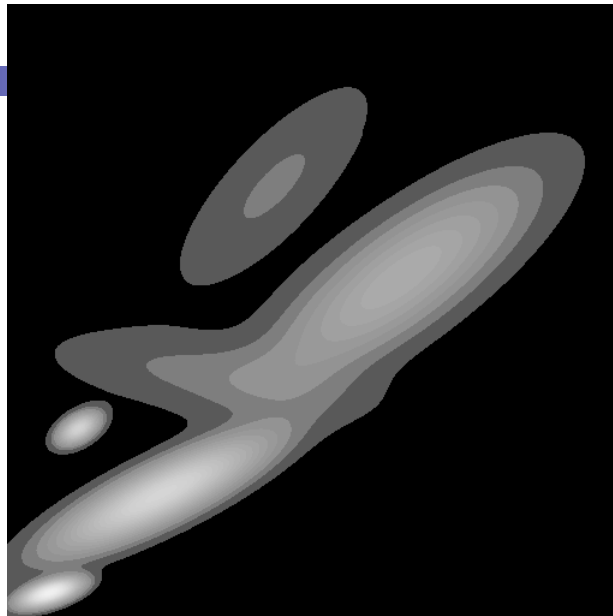
---

# 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*    *← what we have*

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

     $\theta$    parameters

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

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

- Special case:    *← non-invertible, deterministic fcn*

$$x = g(y)$$

   *e.g.* $\quad y = \begin{bmatrix} z \\ x \end{bmatrix}$ *← class labels*    *in standard*

              *← obs.*    *mix. models*

# Expectation Maximization (EM) – Derivation

- Step 1    $p(y, x \mid \theta)$ *← bc $x = g(y)$*
  - Rewrite desired likelihood in terms of complete data terms

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

           *quantity of interest*

$$\Rightarrow \log p(x \mid \theta) = \log p(y \mid \theta) - \log p(y \mid x, \theta)$$

     $\underbrace{\qquad}_{L_x(\theta)}$

- Step 2
  - Assume estimate of parameters $\hat{\theta}$
  - Take expectation with respect to $p(y \mid x, \hat{\theta})$    *"$E[\cdot \mid x, \hat{\theta}]$"*

$$L_x(\theta) = \underbrace{E\left[\log p(y \mid \theta) \mid x, \hat{\theta}\right]}_{U(\theta, \hat{\theta})} + \underbrace{E\left[-\log p(y \mid x, \theta) \mid x, \hat{\theta}\right]}_{V(\theta, \hat{\theta})}$$

# Expectation Maximization (EM) – Derivation

- **Step 3**
  - Consider log likelihood of data at any $\theta$ relative to log likelihood at $\hat{\theta}$

$$L_x(\theta) - L_x(\hat{\theta}) = \left[U(\theta,\hat{\theta}) - U(\hat{\theta},\hat{\theta})\right] + \left[V(\theta,\hat{\theta}) - V(\hat{\theta},\hat{\theta})\right]$$
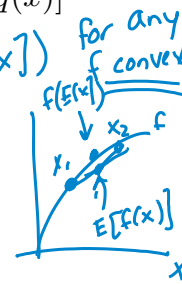
- **Aside: Gibbs Inequality** $E_p[\log p(x)] \geq E_p[\log q(x)]$ ✓

Proof: Use Jensen's Ineq. $E[f(x)] \leq f(E[x])$ for any f convex

$f(E[x])$   $x_2$  f
$x_1$
$E[f(x)]$
$x$

Here:

$$E_p[\log q] - E_p[\log p] = E_p\left[\log \frac{q}{p}\right]$$
$$\leq \log E_p\left[\frac{q}{p}\right] = \log \int_x p(x) \frac{q(x)}{p(x)} dx = 0$$

37

---

# Expectation Maximization (EM) – Derivation

$$L_x(\theta) - L_x(\hat{\theta}) = [U(\theta,\hat{\theta}) - U(\hat{\theta},\hat{\theta})] + \underbrace{[V(\theta,\hat{\theta}) - V(\hat{\theta},\hat{\theta})]}_{\geq 0}$$

- **Step 4**
  - Determine conditions under which log likelihood at $\theta$ exceeds that at $\hat{\theta}$

Using Gibbs inequality:

$$V(\theta,\hat{\theta}) = E\left[-\log p(y|x,\theta) \,|\, x,\hat{\theta}\right] \geq E\left[-\log p(y|x,\hat{\theta}) \,|\, x,\hat{\theta}\right]$$
$$= V(\hat{\theta},\hat{\theta}) \quad \forall \theta$$

If $U(\theta,\hat{\theta}) \geq U(\hat{\theta},\hat{\theta})$

choosing $\theta$ s.t. this is true means we're moving in the right direction (or at least not wrong)

Then

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

38

19

# Motivates EM Algorithm

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

- **E-Step**

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

- **M-Step**

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

  From before, $U(\hat{\theta}^{(t+1)}, \hat{\theta}^{(t)}) \geq U(\hat{\theta}^{(t)}, \hat{\theta}^{(t)})$

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

39

# Example – Mixture Models

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

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

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

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

40

20

# Coordinate Ascent Behavior

- Bound log likelihood:

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

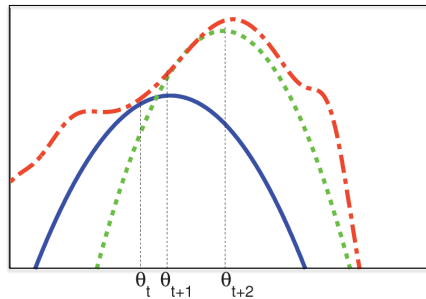

Figure from
KM textbook

---

# 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 $p(y \mid \theta)$ in *U* with $p(y, x \mid \theta)$

# 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 "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

# 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