# EM

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University
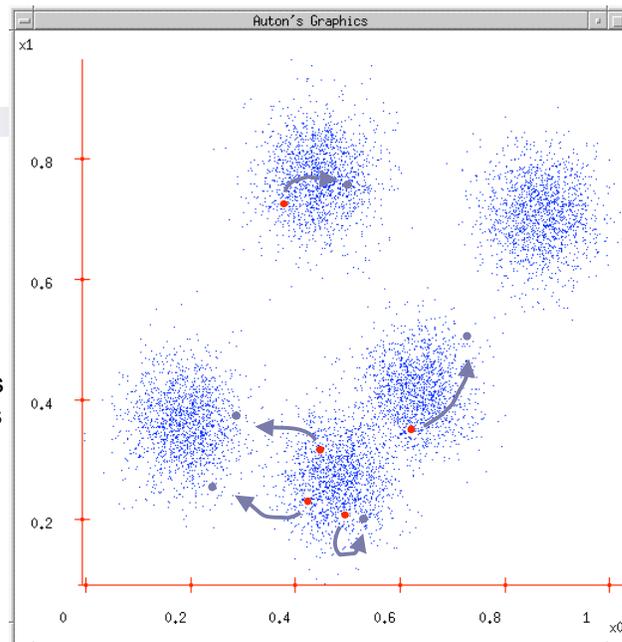
November 19th, 2007

1

---

# 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…

5. …and jumps there

6. …Repeat until terminated!

2

# K-means

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

- **Classify**: Assign each point j∈{1,…m} to nearest center: *center of point j is closest to j*
  - $C^{(t)}(j) \leftarrow \arg\min_i ||\mu_i - x_j||^2$

- **Recenter**: $\mu_i$ becomes centroid of its point:
  - $\mu_i^{(t+1)} \leftarrow \arg\min_\mu \sum_{j:C(j)=i} ||\mu - x_j||^2$    *opt $\mu_i = \sum_{j:c(j)=i} x_j$ is the mean!!*
  $$\frac{\sum_{j:c(j)=i} x_j}{\sum_{j:c(s)=i} 1}$$
  - Equivalent to $\mu_i \leftarrow$ average of its points! ☺

3

---

# Does K-means converge??? Part 2

- Optimize potential function:
$$\min_\mu \min_C F(\mu, C) = \min_\mu \min_C \sum_{i=1}^{k} \sum_{j:C(j)=i} ||\mu_i - x_j||^2$$
  $C = \bar{C}$
- Fix C, optimize $\mu$

$$\min_\mu \sum_{i=1}^{K} \sum_{j:\bar{C}(j)=i} ||\mu_i - x_j||^2$$

$$= \sum_{i=1}^{K} \min_{\mu_i} \sum_{j:\bar{C}(j)=i} ||\mu_i - x_j||^2$$
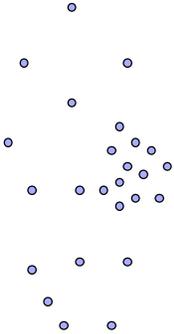
*$\mu_i$ is mean of points in cluster i*   ⟹   *re center in K-means.*

4

# Coordinate descent algorithms

$$\min_{\mu} \min_{C} F(\mu, C) = \min_{\mu} \min_{C} \sum_{i=1}^{k} \sum_{j:C(j)=i} ||\mu_i - x_j||^2$$

- Want: $\min_a \min_b F(a,b)$
- Coordinate descent:
  - fix a, minimize b
  - fix b, minimize a
  - repeat
- Converges!!!
  - if F is bounded
  - to a (often good) local optimum
    - as we saw in applet (play with it!)

- K-means is a coordinate descent algorithm!

# (One) bad case for k-means

- Clusters may overlap
- Some clusters may be "wider" than others
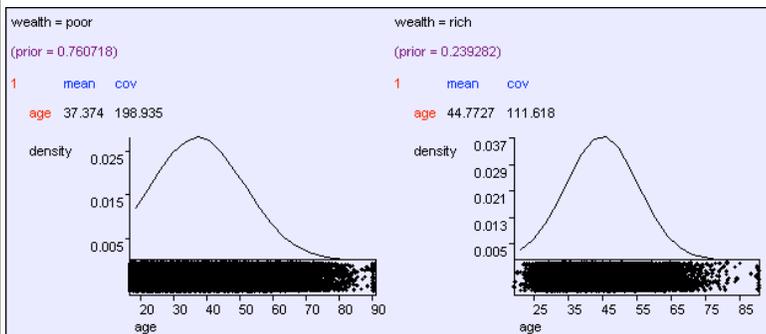
# Gaussian Bayes Classifier Reminder

$$P(y = i \mid \mathbf{x}_j) = \frac{p(\mathbf{x}_j \mid y = i)P(y = i)}{p(\mathbf{x}_j)}$$

$$P(y = i \mid \mathbf{x}_j) \propto \frac{1}{(2\pi)^{m/2} \parallel \Sigma_i \parallel^{1/2}} \exp\left[-\frac{1}{2}\left(\mathbf{x}_j - \mu_i\right)^T \Sigma_i^{-1}\left(\mathbf{x}_j - \mu_i\right)\right] P(y = i)$$
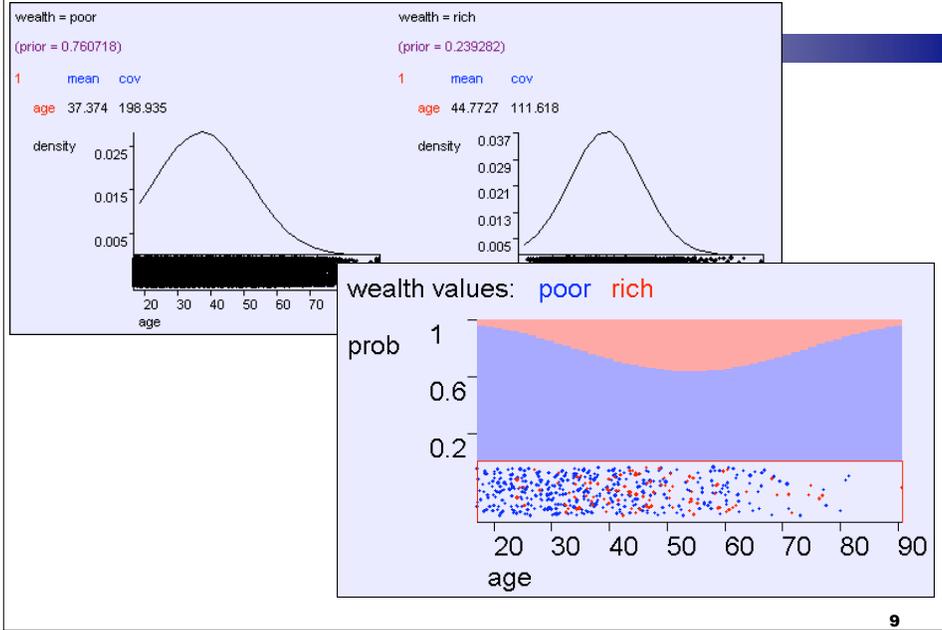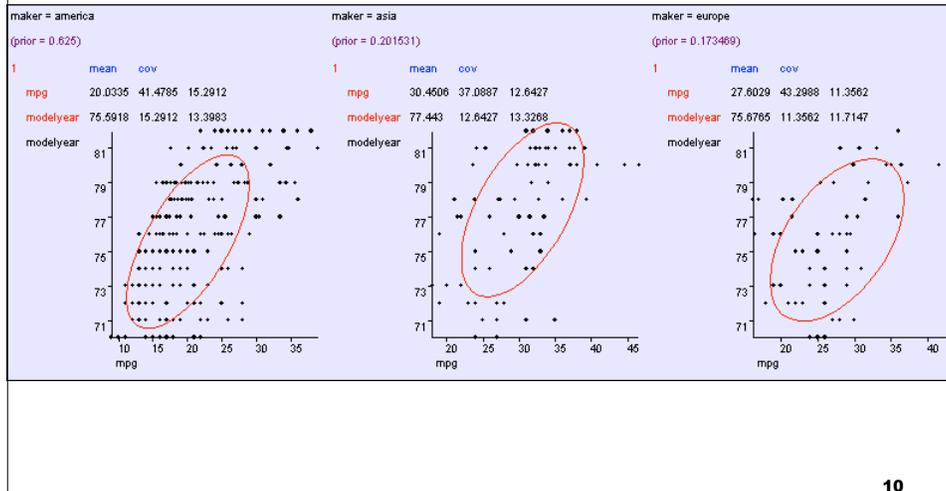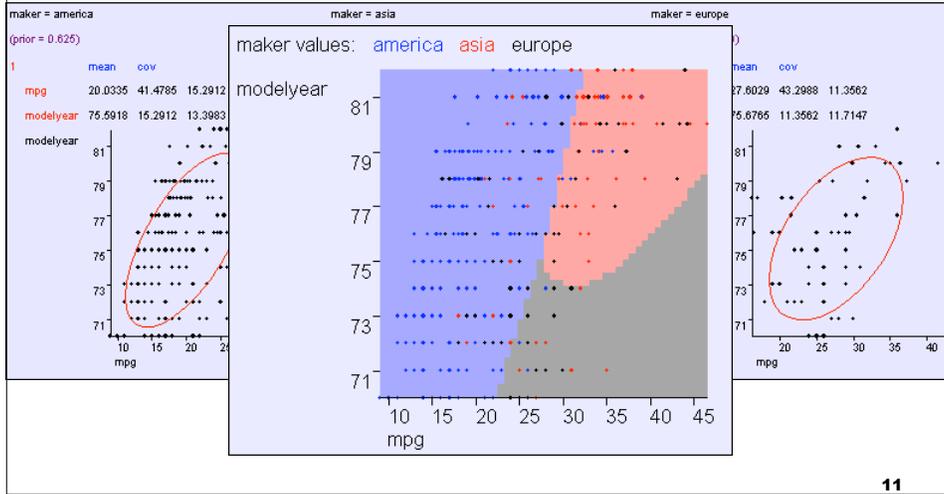
# Predicting wealth from age

# Predicting wealth from age

# Learning modelyear , mpg ---> maker

$$\Sigma = \begin{pmatrix} \sigma^2_1 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma^2_2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma^2_m \end{pmatrix}$$
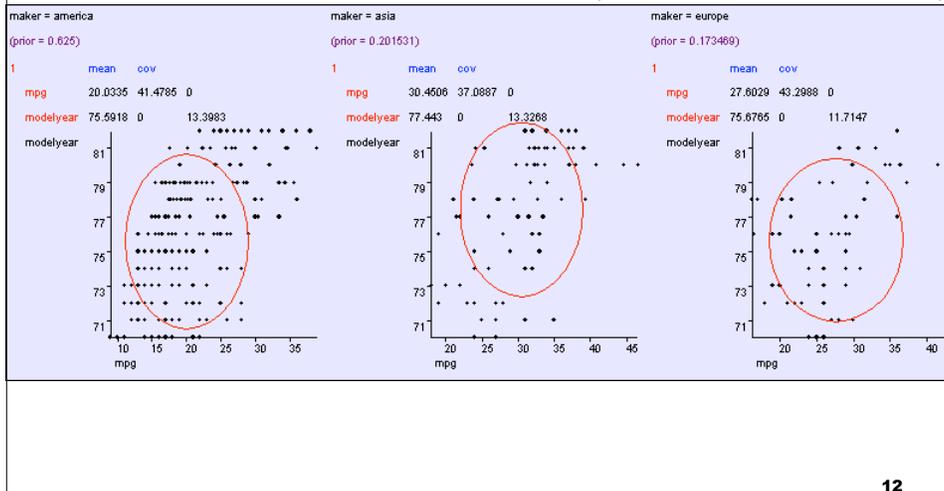
# General: $O(m^2)$ parameters

$$\Sigma = \begin{pmatrix} \sigma^2_1 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma^2_2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma^2_m \end{pmatrix}$$

# Aligned: $O(m)$ parameters

$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma^2_m \end{pmatrix}$$

# Aligned: *O(m)* parameters

$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma^2_m \end{pmatrix}$$

13



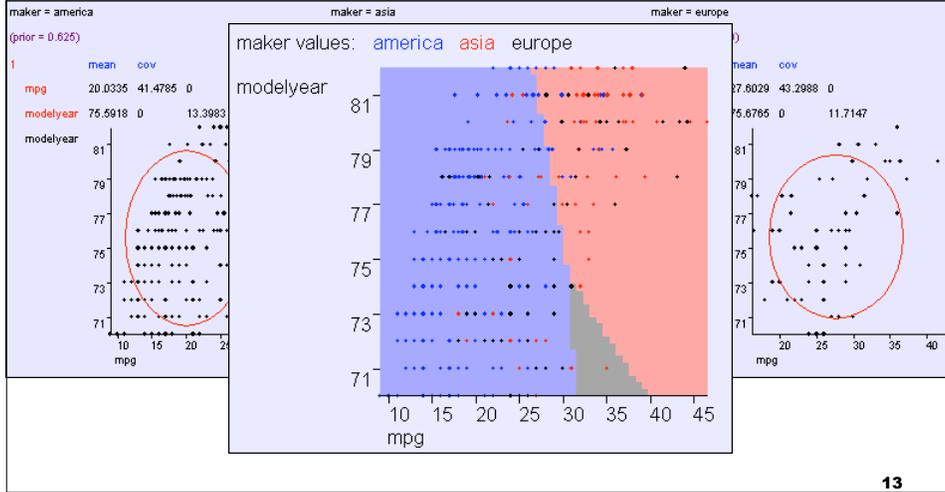# Spherical: *O(1)* cov parameters

$$\Sigma = \begin{pmatrix} \sigma^2 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma^2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma^2 \end{pmatrix}$$
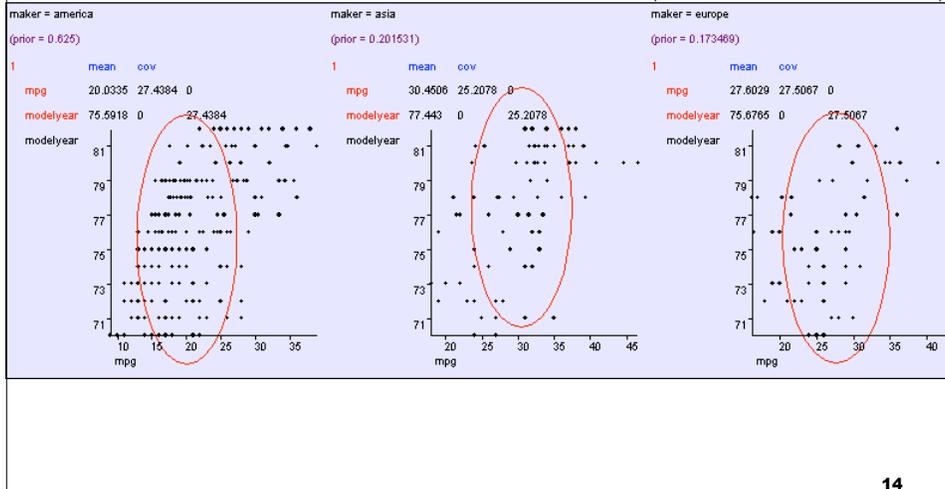
14

# Spherical: *O(1)* cov parameters

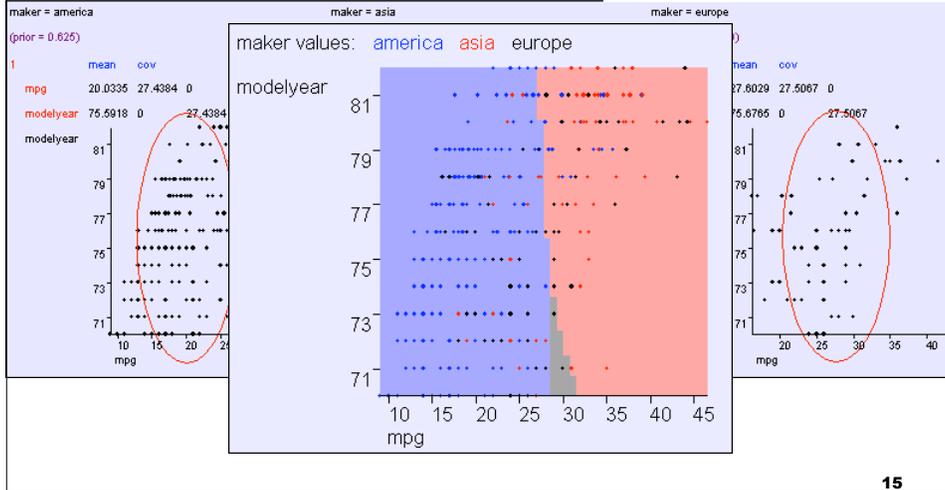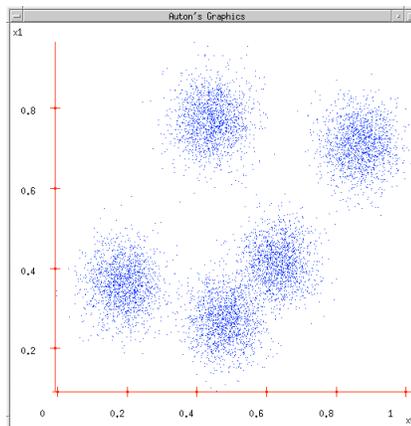$$\Sigma = \begin{pmatrix} \sigma^2 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma^2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma^2 \end{pmatrix}$$



15

---

## Next…   back to Density Estimation

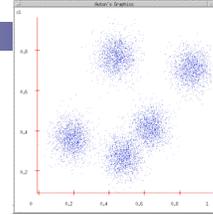What if we want to do density estimation with multimodal or clumpy data?



16

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



- MLE:
  - argmax $\prod_j P(y_j, x_j)$

- But we don't know $y_j$'s!!!
- Maximize marginal likelihood:
  - argmax $\prod_j P(x_j)$ = argmax $\prod_j \sum_{i=1}^{k} P(y_j=i, x_j)$

# Special case: spherical Gaussians and hard assignments

$$P(y = i \mid \mathbf{x}_j) \propto \frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}\left(\mathbf{x}_j - \mu_i\right)^T \Sigma_i^{-1}\left(\mathbf{x}_j - \mu_i\right)\right] P(y = i)$$

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

$$P(\mathbf{x}_j \mid y = i) \propto \exp\left[-\frac{1}{2\sigma^2}\left\|\mathbf{x}_j - \mu_i\right\|^2\right]$$

- If each $x_j$ belongs to one class $C(j)$ (hard assignment), marginal likelihood:

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

- Same as K-means!!!

# The GMM assumption

- There are k components

- Component *i* has an associated mean vector $\mu_i$

$$\mu_2$$
$$\mu_1$$
$$\mu_3$$

19

---

# The GMM assumption

• There are k components

• Component *i* has an associated mean vector $\mu_i$

• Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$

  Each data point is generated according to the following recipe:

$$\mu_2$$
$$\mu_1$$
$$\mu_3$$

20

# The GMM assumption

- There are k components

- Component *i* has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 \boldsymbol{I}$

Each data point is generated according to the following recipe:

1. Pick a component at random: Choose component i with probability $P(y=i)$

$\mu_2$

# The GMM assumption

- There are k components

- Component *i* has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 \boldsymbol{I}$

Each data point is generated according to the following recipe:

1. Pick a component at random: Choose component i with probability $P(y=i)$

2. Datapoint ~ $N(\mu_i, \sigma^2 \boldsymbol{I})$
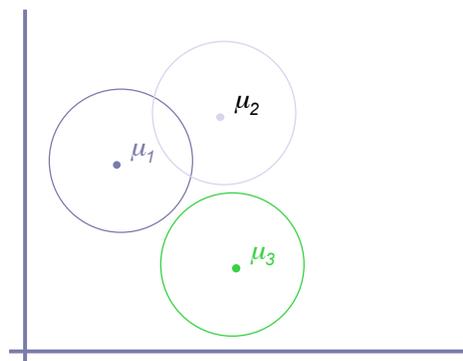
$\mu_2$

x

# The General GMM assumption

- There are k components

- Component $i$ has an associated mean vector $\mu_i$

- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\Sigma_i$

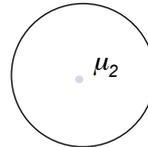Each data point is generated according to the following recipe:

1. Pick a component at random: Choose component i with probability $P(y=i)$
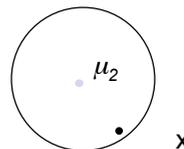
2. Datapoint ~ $N(\mu_i, \Sigma_i)$

$\mu_2$

$\mu_1$

$\mu_3$

23

---

# Unsupervised Learning: not as hard as it looks

Sometimes easy

Sometimes impossible

and sometimes in between

*IN CASE YOU'RE WONDERING WHAT THESE DIAGRAMS ARE, THEY SHOW 2-d UNLABELED DATA (X VECTORS) DISTRIBUTED IN 2-d SPACE. THE TOP ONE HAS THREE VERY CLEAR GAUSSIAN CENTERS*

24

# Marginal likelihood for general case

$$P(y = i \mid \mathbf{x}_j) \propto \frac{1}{(2\pi)^{m/2} \| \Sigma_i \|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}_j - \mu_i)\right] P(y = i)$$

- Marginal likelihood:

$$\prod_{j=1}^m P(\mathbf{x}_j) = \prod_{j=1}^m \sum_{i=1}^k P(\mathbf{x}_j, y = i)$$

$$= \prod_{j=1}^m \sum_{i=1}^k \frac{1}{(2\pi)^{m/2} \| \Sigma_i \|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}_j - \mu_i)\right] P(y = i)$$

# Special case 2: spherical Gaussians and soft assignments

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

$$P(\mathbf{x}_j \mid y = i) \propto \exp\left[-\frac{1}{2\sigma^2}\|\mathbf{x}_j - \mu_i\|^2\right]$$

- Uncertain about class of each $x_j$ (soft assignment), marginal likelihood:

$$\prod_{j=1}^m \sum_{i=1}^k P(\mathbf{x}_j, y = i) \propto \prod_{j=1}^m \sum_{i=1}^k \exp\left[-\frac{1}{2\sigma^2}\|\mathbf{x}_j - \mu_i\|^2\right] P(y = i)$$

# Unsupervised Learning: Mediumly Good News

We now have a procedure s.t. if you give me a guess at $\mu_1, \mu_2 .. \mu_k,$

I can tell you the prob of the unlabeled data given those $\mu$'s.

Suppose $x$'s are 1-dimensional.

**(From Duda and Hart)**

There are two classes; $w_1$ and $w_2$

$P(y_1) = 1/3 \quad P(y_2) = 2/3 \quad \sigma = 1$ .

There are 25 unlabeled datapoints

$x_1 = 0.608$
$x_2 = -1.590$
$x_3 = 0.235$
$x_4 = 3.949$
      :
$x_{25} = -0.712$



DATA SCATTERGRAM

-4    -2    0    2    4

---

# Duda & Hart's Example

We can graph the prob. dist. function of data given our $\mu_1$ and $\mu_2$ estimates.

We can also graph the true function from which the data was randomly generated.



- They are close. Good.

- The 2nd solution tries to put the "2/3" hump where the "1/3" hump should go, and vice versa.

- In this example unsupervised is almost as good as supervised. If the $x_1 .. x_{25}$ are given the class which was used to learn them, then the results are ($\mu_1 = -2.176, \mu_2 = 1.684$). Unsupervised got ($\mu_1 = -2.13, \mu_2 = 1.668$).

# Duda & Hart's Example



Graph of
log P($x_1$, $x_2$ .. $x_{25}$ | $\mu_1$, $\mu_2$ )
   against $\mu_1$ ($\rightarrow$) and $\mu_2$ ($\uparrow$)

Max likelihood = ($\mu_1$ =-2.13, $\mu_2$ =1.668)

Local minimum, but very close to global at ($\mu_1$ =2.085, $\mu_2$ =-1.257)*

    * corresponds to switching $y_1$ with $y_2$.

---

# Finding the max likelihood $\mu_1, \mu_2 .. \mu_k$

We can compute  P( data | $\mu_1, \mu_2 .. \mu_k$)

How do we find the $\mu_i$'s which give max. likelihood?

- The normal max likelihood trick:
  Set $\dfrac{\partial}{\partial \mu_i}$ log Prob (….) = 0
  and solve for $\mu_i$'s.
    - \# Here you get non-linear non-analytically-solvable equations
- Use gradient descent
  Often slow but doable
- Use a much faster, cuter, and recently very popular method…

# Announcements

- HW5 out later today…
  - Due December 5th by 3pm to Monica Hopes, Wean 4619
- Project:
  - Poster session: NSH Atrium, Friday 11/30, 2-5pm
    - Print your poster early!!!
      - SCS facilities has a poster printer, ask helpdesk
      - Students from outside SCS should check with their departments
      - It's OK to print separate pages
    - We'll provide pins, posterboard and an easel
      - Poster size: 32x40 inches
    - Invite your friends, there will be a prize for best poster, by popular vote
- Last lecture:
  - Thursday, 11/29, 5-6:20pm, Wean 7500

---



Expectation
Maximalization

# The E.M. Algorithm

DETOUR

- We'll get back to unsupervised learning soon
- But now we'll look at an even simpler case with hidden information
- The EM algorithm
  - Can do trivial things, such as the contents of the next few slides
  - An excellent way of doing our unsupervised learning problem, as we'll see
  - Many, many other uses, including learning BNs with hidden data

# Silly Example

Let events be "grades in a class"

| | |
|---|---|
| $w_1$ = Gets an A | $P(A) = \frac{1}{2}$ |
| $w_2$ = Gets a  B | $P(B) = \mu$ |
| $w_3$ = Gets a  C | $P(C) = 2\mu$ |
| $w_4$ = Gets a  D | $P(D) = \frac{1}{2} - 3\mu$ |

(Note $0 \le \mu \le 1/6$)

Assume we want to estimate $\mu$ from data.  In a given class there were

a   A's
b   B's
c   C's
d   D's

What's the maximum likelihood estimate of $\mu$ given a,b,c,d ?

# Trivial Statistics

P(A) = ½   P(B) = μ   P(C) = 2μ   P(D) = ½-3μ

P( $a,b,c,d$ | μ) = K(½)$^a$(μ)$^b$(2μ)$^c$(½-3μ)$^d$

log P( $a,b,c,d$ | μ) = log K + $a$log ½ + $b$log μ + $c$log 2μ + $d$log (½-3μ)

FOR MAX LIKE $\mu$, SET $\dfrac{\partial \text{LogP}}{\partial \mu} = 0$

$$\frac{\partial \text{LogP}}{\partial \mu} = \frac{b}{\mu} + \frac{2c}{2\mu} - \frac{3d}{1/2 - 3\mu} = 0$$

Gives max like $\mu = \dfrac{b+c}{6(b+c+d)}$

So if class got

| A | B | C | D |
|---|---|---|---|
| 14 | 6 | 9 | 10 |

Max like $\mu = \dfrac{1}{10}$

Boring, but true!

---

# Same Problem with Hidden Information

Someone tells us that

Number of High grades (A's + B's) = $h$

Number of C's            = $c$

Number of D's            = $d$

What is the max. like estimate of μ now?

REMEMBER
P(A) = ½
P(B) = μ
P(C) = 2μ
P(D) = ½-3μ

# Same Problem with Hidden Information

Someone tells us that

Number of High grades (A's + B's) = $h$

Number of C's = $c$

Number of D's = $d$

What is the max. like estimate of µ now?

We can answer this question circularly:

**EXPECTATION**  If we know the value of µ we could compute the expected value of $a$ and $b$

Since the ratio a:b should be the same as the ratio ½ : µ

$$a = \frac{\frac{1}{2}}{\frac{1}{2} + \mu} h \qquad b = \frac{\mu}{\frac{1}{2} + \mu} h$$

**MAXIMIZATION**

If we know the expected values of $a$ and $b$ we could compute the maximum likelihood value of µ

$$\mu = \frac{b + c}{6(b + c + d)}$$

37

---

# E.M. for our Trivial Problem

We begin with a guess for µ
We iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of µ and $a$ and $b$.

Define  $\mu^{(t)}$ the estimate of µ on the t'th iteration
 $b^{(t)}$ the estimate of $b$ on t'th iteration

$$\mu^{(0)} = \text{initial guess}$$

$$b^{(t)} = \frac{\mu^{(t)} h}{\frac{1}{2} + \mu^{(t)}} = \mathrm{E}\left[b \mid \mu^{(t)}\right]$$  **E-step**

$$\mu^{(t+1)} = \frac{b^{(t)} + c}{6(b^{(t)} + c + d)}$$  **M-step**

$$= \text{max like est. of } \mu \text{ given } b^{(t)}$$

**Continue iterating until converged.**
**Good news: Converging to local optimum is assured.**
**Bad news: I said "local" optimum.**
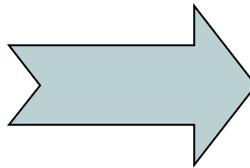
38

# E.M. Convergence

- Convergence proof based on fact that Prob(data | μ) must increase or remain same between each iteration [NOT OBVIOUS]
- But it can never exceed 1   [OBVIOUS]

So it must therefore converge   [OBVIOUS]

In our example, suppose we had

$$h = 20$$
$$c = 10$$
$$d = 10$$
$$\mu^{(0)} = 0$$

Convergence is generally <u>linear</u>: error decreases by a constant factor each time step.

| t | $\mu^{(t)}$ | $b^{(t)}$ |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 0.0833 | 2.857 |
| 2 | 0.0937 | 3.158 |
| 3 | 0.0947 | 3.185 |
| 4 | 0.0948 | 3.187 |
| 5 | 0.0948 | 3.187 |
| 6 | 0.0948 | 3.187 |

# Back to Unsupervised Learning of GMMs – a simple case

A simple case:

We have unlabeled data $x_1\ x_2\ \dots\ x_m$
We know there are k classes
We know $P(y_1)\ P(y_2)\ P(y_3)\ \dots\ P(y_k)$
We <u>don't</u> know $\mu_1\ \mu_2\ ..\ \mu_k$

We can write P( data | $\mu_1 \dots \mu_k$ )

$$= p\left(x_1 \dots x_m \middle| \mu_1 \dots \mu_k\right)$$

$$= \prod_{j=1}^{m} p\left(x_j \middle| \mu_1 \dots \mu_k\right)$$

$$= \prod_{j=1}^{m}\sum_{i=1}^{k} p\left(x_j \middle| \mu_i\right)P(y = i)$$

$$\propto \prod_{j=1}^{m}\sum_{i=1}^{k} \exp\left(-\frac{1}{2\sigma^2}\left\|x_j - \mu_i\right\|^2\right)P(y = i)$$

# EM for simple case of GMMs: The E-step

- If we know $\mu_1, \ldots, \mu_k$ $\rightarrow$ easily compute prob. point $x_j$ belongs to class $y=i$

$$p\left(y = i \middle| x_j, \mu_1 \ldots \mu_k\right) \propto \exp\left(-\frac{1}{2\sigma^2}\left\|x_j - \mu_i\right\|^2\right) P(y = i)$$

# EM for simple case of GMMs: The M-step

- If we know prob. point $x_j$ belongs to class $y=i$
  $\rightarrow$ MLE for $\mu_i$ is weighted average
  - imagine k copies of each $x_j$, each with weight $P(y=i|x_j)$:

$$\mu_i = \frac{\sum_{j=1}^{m} P\left(y = i \middle| x_j\right) x_j}{\sum_{j=1}^{m} P\left(y = i \middle| x_j\right)}$$

# E.M. for GMMs

**E-step**

Compute "expected" classes of all datapoints for each class

$$p\left(y = i \mid x_j, \mu_1 \ldots \mu_k\right) \propto \exp\left(-\frac{1}{2\sigma^2} \left\| x_j - \mu_i \right\|^2\right) P(y = i)$$

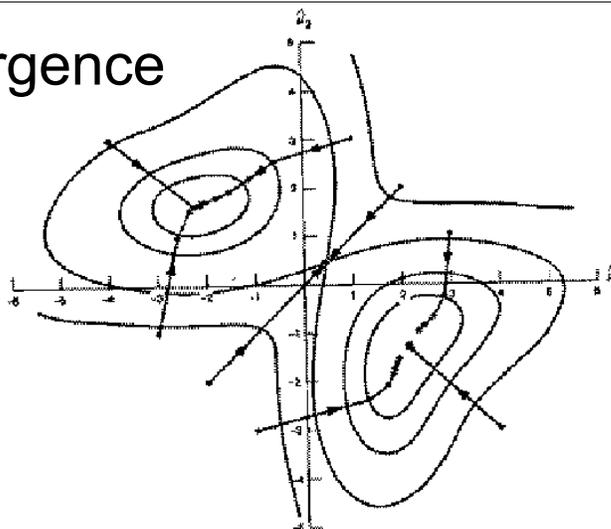*Just evaluate a Gaussian at $x_j$*

**M-step**

Compute Max. like **μ** given our data's class membership distributions

$$\mu_i = \frac{\sum_{j=1}^{m} P\left(y = i \mid x_j\right) x_j}{\sum_{j=1}^{m} P\left(y = i \mid x_j\right)}$$

43

---

# E.M. Convergence



- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func. ! convergence to a local optimum guaranteed
- See Neal & Hinton reading on class webpage

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

44

# E.M. for axis-aligned GMM

$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma^2_m \end{pmatrix}$$

Iterate. On the $t$'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \dots p_k^{(t)} \}$$

$p_i^{(t)}$ is shorthand for estimate of $P(y=i)$ on t'th iteration

**E-step**

Compute "expected" classes of all datapoints for each class

$$P\left(y = i \middle| x_j, \lambda_t \right) \propto p_i^{(t)} p\left(x_j \middle| \mu_i^{(t)}, \Sigma_i^{(t)} \right)$$

*Just evaluate a Gaussian at $x_j$*

M-step

Compute Max. like **μ** given our data's class membership distributions

$$\grave{\imath}_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x_j, \lambda_t \right) x_j}{\sum_j P\left(y = i \middle| x_j, \lambda_t \right)}$$

$$p_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x_j, \lambda_t \right)}{m}$$

$m$ = #records

45

---

# E.M. for General GMMs

$p_i^{(t)}$ is shorthand for estimate of $P(y=i)$ on t'th iteration

Iterate. On the $t$'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \dots p_k^{(t)} \}$$

**E-step**

Compute "expected" classes of all datapoints for each class

$$P\left(y = i \middle| x_j, \lambda_t \right) \propto p_i^{(t)} p\left(x_j \middle| \mu_i^{(t)}, \Sigma_i^{(t)} \right)$$

*Just evaluate a Gaussian at $x_j$*

M-step

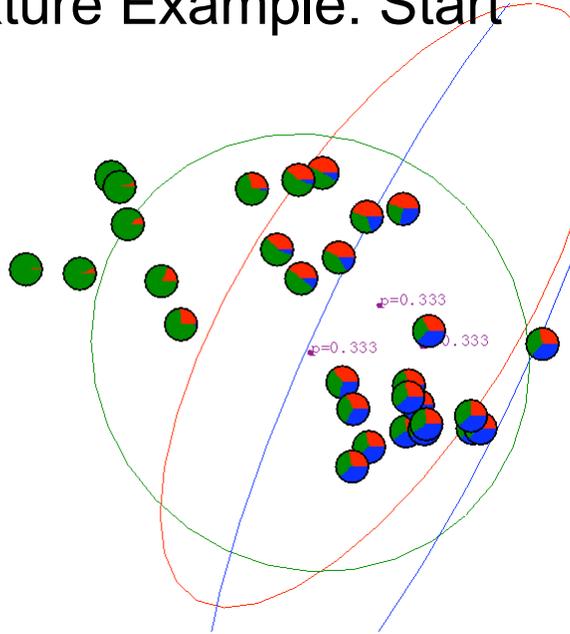Compute Max. like **μ** given our data's class membership distributions

$$\grave{\imath}_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x_j, \lambda_t \right) x_j}{\sum_j P\left(y = i \middle| x_j, \lambda_t \right)} \qquad \Sigma_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x_j, \lambda_t \right)\left[x_j - \mu_i^{(t+1)}\right]\left[x_j - \mu_i^{(t+1)}\right]}{\sum_j P\left(y = i \middle| x_j, \lambda_t \right)}$$

$$p_i^{(t+1)} = \frac{\sum_j P\left(y = i \middle| x_j, \lambda_t \right)}{m}$$

$m$ = #records
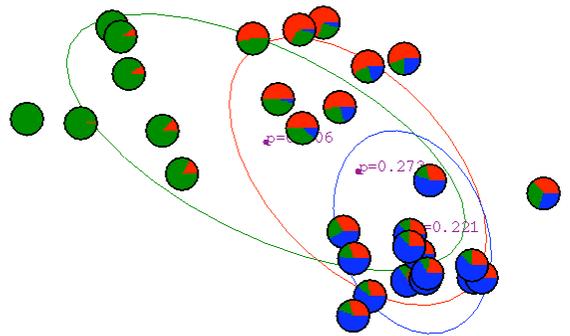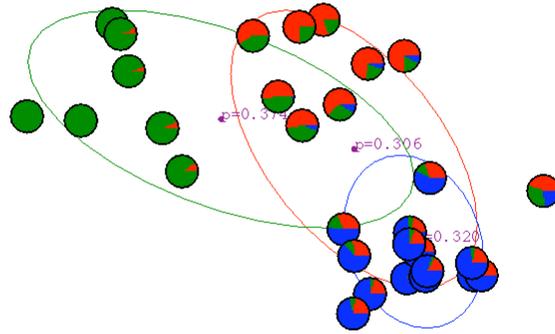
46

Gaussian Mixture Example: Start
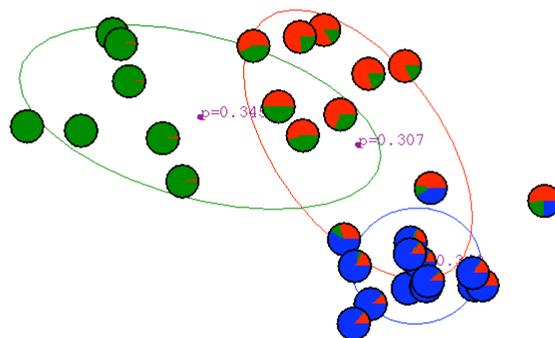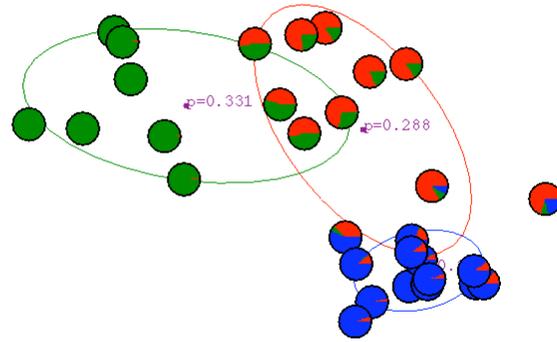

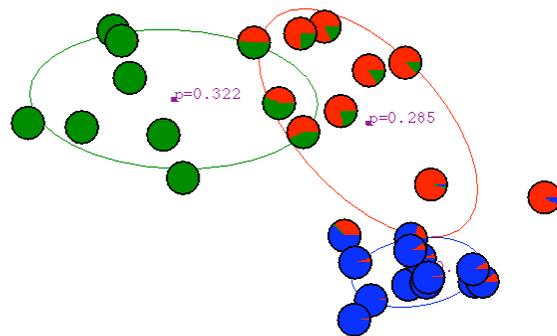After first iteration

# After 2nd iteration



49
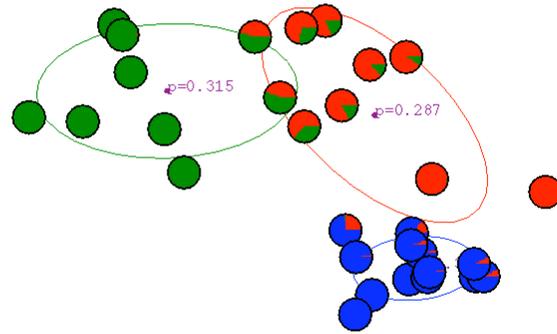
# After 3rd iteration



50

# After 4th iteration



p=0.331

p=0.288

51

# After 5th iteration



p=0.322

p=0.285

52

# After 6th iteration



p=0.315

p=0.287

# After 20th iteration



p=0.234

p=0.334

# Some Bio Assay data

# GMM clustering of the assay data

# Resulting Density Estimator



57

# Three classes of assay

(each learned with it's own mixture model)



Compound =
IL-1
TNF
none

58

# Resulting Bayes Classifier

Compound =
IL-1
TNF
none

nucleus

59

---

# Resulting Bayes Classifier, using posterior probabilities to alert about ambiguity and anomalousness

Yellow means anomalous

Cyan means ambiguous

Compound =
IL-1
TNF
none

Yellow means ANOMALOUS
Cyan means AMBIGUOUS

nucleus

60

# The general learning problem with missing data

- Marginal likelihood – **x** is observed, **z** is missing:

$$\ell(\theta : \mathcal{D}) = \log \prod_{j=1}^{m} P(\mathbf{x}_j \mid \theta)$$

$$= \sum_{j=1}^{m} \log P(\mathbf{x}_j \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)$$

# E-step

- **x** is observed, **z** is missing
- Compute probability of missing data given current choice of $\theta$
  - $Q(\mathbf{z}|\mathbf{x}_j)$ for each $\mathbf{x}_j$
    - e.g., probability computed during classification step
    - corresponds to "classification step" in K-means

$$Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) = P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})$$

# Jensen's inequality

$$\ell(\theta : \mathcal{D}) \;=\; \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{z} \mid \mathbf{x}_j) P(\mathbf{x}_j \mid \theta)$$

- **Theorem**: $\log \sum_{\mathbf{z}} P(\mathbf{z})\, f(\mathbf{z}) \;\geq\; \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

# Applying Jensen's inequality

- Use:  $\log \sum_{\mathbf{z}} P(\mathbf{z})\, f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\ell(\theta^{(t)} : \mathcal{D}) \;=\; \sum_{j=1}^{m} \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j)}$$

# The M-step maximizes lower bound on weighted data

- Lower bound from Jensen's:

$$\ell(\theta^{(t)} : \mathcal{D}) \;\geq\; \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Corresponds to weighted dataset:
  - $<\mathbf{x}_1, \mathbf{z}=1>$ with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}_1)$
  - $<\mathbf{x}_1, \mathbf{z}=2>$ with weight $Q^{(t+1)}(\mathbf{z}=2|\mathbf{x}_1)$
  - $<\mathbf{x}_1, \mathbf{z}=3>$ with weight $Q^{(t+1)}(\mathbf{z}=3|\mathbf{x}_1)$
  - $<\mathbf{x}_2, \mathbf{z}=1>$ with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}_2)$
  - $<\mathbf{x}_2, \mathbf{z}=2>$ with weight $Q^{(t+1)}(\mathbf{z}=2|\mathbf{x}_2)$
  - $<\mathbf{x}_2, \mathbf{z}=3>$ with weight $Q^{(t+1)}(\mathbf{z}=3|\mathbf{x}_2)$
  - …

# The M-step

$$\ell(\theta^{(t)} : \mathcal{D}) \;\geq\; \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta)$$

- Use expected counts instead of counts:
  - If learning requires Count($\mathbf{x}$,$\mathbf{z}$)
  - Use $E_{Q(t+1)}$[Count($\mathbf{x}$,$\mathbf{z}$)]

# Convergence of EM

- Define potential function $F(\theta, Q)$:

$$\ell(\theta : \mathcal{D}) \;\geq\; F(\theta, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

- EM corresponds to coordinate ascent on F
  - Thus, maximizes lower bound on marginal log likelihood

# M-step is easy

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta)$$

- Using potential function

$$F(\theta, Q^{(t+1)}) \;=\; \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t+1)})$$

# E-step also doesn't decrease potential function 1

- Fixing θ to θ^(t):

$$\ell(\theta^{(t)} : \mathcal{D}) \;\geq\; F(\theta^{(t)}, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

# KL-divergence

- Measures distance between distributions

$$KL(Q||P) = \sum_{z} Q(z) \log \frac{Q(z)}{P(z)}$$

- KL=zero if and only if Q=P

# E-step also doesn't decrease potential function 2

- Fixing θ to θ(t):

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) \;=\; \ell(\theta^{(t)} : \mathcal{D}) + \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

$$= \; \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^{m} KL\left(Q(\mathbf{z} \mid \mathbf{x}_j) || P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})\right)$$

---

# E-step also doesn't decrease potential function 3

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) \;=\; \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^{m} KL\left(Q(\mathbf{z} \mid \mathbf{x}_j) || P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})\right)$$

- Fixing θ to θ(t)
- Maximizing F(θ(t),Q) over Q → set Q to posterior probability:

$$Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \;\leftarrow\; P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})$$

- Note that

$$F(\theta^{(t)}, Q^{(t+1)}) \;=\; \ell(\theta^{(t)} : \mathcal{D})$$

# EM is coordinate ascent

$$\ell(\theta : \mathcal{D}) \ \geq \ F(\theta, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

- **M-step**: Fix Q, maximize F over θ (a lower bound on $\ell(\theta : \mathcal{D})$ ):

$$\ell(\theta : \mathcal{D}) \ \geq \ F(\theta, Q^{(t)}) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t)})$$

- **E-step**: Fix θ, maximize F over Q:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) \ = \ \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^{m} KL\left(Q(\mathbf{z} \mid \mathbf{x}_j) \| P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})\right)$$

  - □ "Realigns" F with likelihood:

$$F(\theta^{(t)}, Q^{(t+1)}) \ = \ \ell(\theta^{(t)} : \mathcal{D})$$

# 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 <u>DOES</u>
- EM is coordinate ascent
- General case for EM

# Acknowledgements

- K-means & Gaussian mixture models presentation contains material from excellent tutorial by Andrew Moore:
  - http://www.autonlab.org/tutorials/
- K-means Applet:
  - http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/AppletKM.html
- Gaussian mixture models Applet:
  - http://www.neurosci.aist.go.jp/%7Eakaho/MixtureEM.html