# Clustering with Gaussian Mixtures

**Ronald J. Williams**
**CSG220**
**Spring 2007**

**Adapted from the Andrew Moore**
**tutorial of the same name**

Nov 10th, 2001

---

# Unsupervised Learning

- You walk into a bar.
  A stranger approaches and tells you:
  "I've got data from k classes. Each class produces observations with a normal distribution and variance $\sigma^2 I$ . Standard simple multivariate gaussian assumptions. I can tell you the probabilities of each class ."

- So far, looks straightforward.
  "I need a maximum likelihood estimate of the $\mu_i$'s ."

- "No problem," you think.
  "There's just one thing. None of the data are labeled. I have datapoints, but I don't know what class they're from (any of them!)

- Uh oh!!

Clustering with Gaussian Mixtures: Slide 2

**1**

# Multivariate Gaussian Density

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} \| \mathbf{\Sigma} \|^{1/2}} \exp\left[ -\frac{1}{2}(\mathbf{x} - \mathbf{\mu})^T \mathbf{\Sigma}(\mathbf{x} - \mathbf{\mu}) \right]$$

where

$\mathbf{\mu} = \text{mean } (m\text{-dimensiona l vector})$

$\mathbf{\Sigma} = \text{covariance } (m \times m \text{ matrix})$

# Predicting wealth from age



1-dimensional Gaussians

**2**

General: $O(m^2)$ parameters

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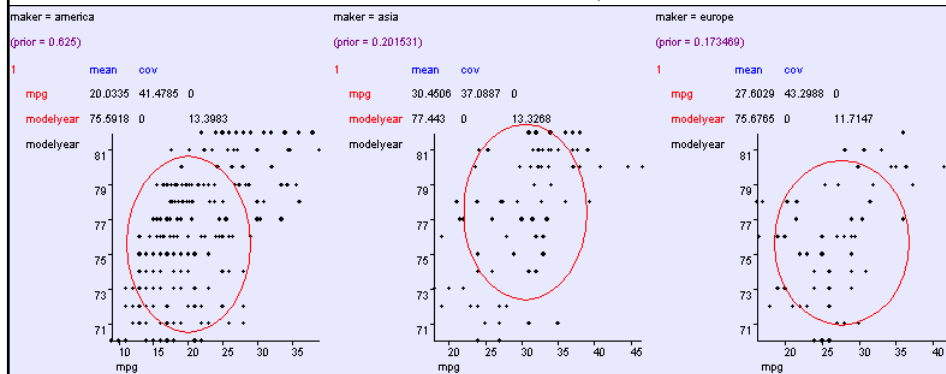
2-dimensional Gaussians

Clustering with Gaussian Mixtures: Slide 5



Aligned: $O(m)$ parameters

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

Clustering with Gaussian Mixtures: Slide 6

# Spherical: *O(1)* cov parameters

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

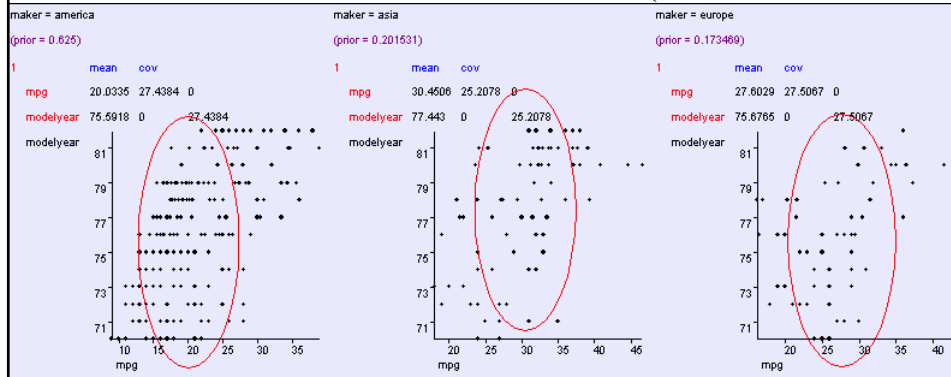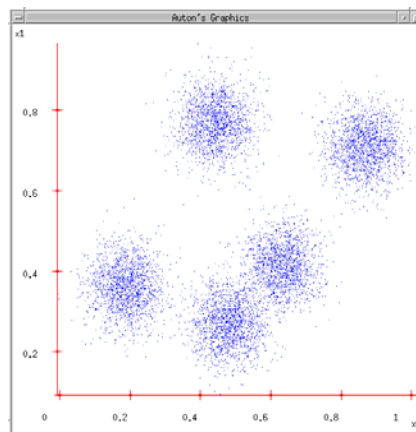# What if we want to do density estimation with <u>multi</u>modal or clumpy data?

Clearly not modeled well by a single Gaussian

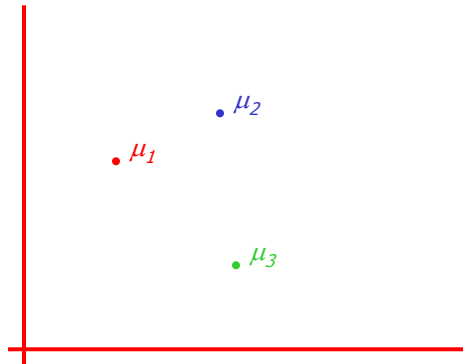# The Gaussian Mixture Model assumption

- There are k components. The i'th component is called $\omega_i$

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

# The Gaussian Mixture Model assumption

- There are k components. The i'th component is called $\omega_i$

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

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

Assume that each datapoint is generated according to the following recipe:

# The Gaussian Mixture Model assumption

- There are k components. The i'th component is called $\omega_i$

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

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

Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.



$\mu_2$

---

# The Gaussian Mixture Model assumption

- There are k components. The i'th component is called $\omega_i$
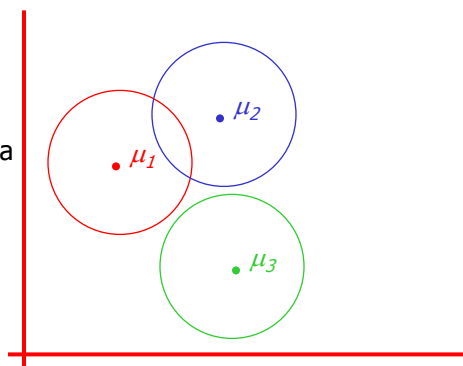
- Component $\omega_i$ has an associated mean vector $\mu_i$

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

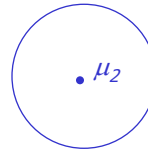Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.

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



$\mu_2$

x

Denotes Gaussian with given mean and covariance
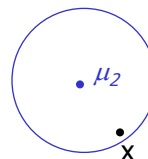
# The General GMM assumption

- There are k components. The i'th component is called $\omega_i$

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

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

Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.

2. Datapoint ~ N($\mu_i$, $\Sigma_i$ )

# 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

## Computing likelihoods in unsupervised case

We have $x_1, x_{2,...} x_R$

We know $P(\omega_1) P(\omega_2) .. P(\omega_k)$

We know $\sigma$

$p(x | \omega_i, \mu_i, ... \mu_k)$ = Prob density that an observation
from class $\omega_i$ would have value $x$
given class means $\mu_1... \mu_k$

Can we write an expression for that?

Clustering with Gaussian Mixtures: Slide 15

---

## Computing likelihoods in unsupervised case

We have $x_1, x_{2,...} x_R$

We know $P(\omega_1) P(\omega_2) .. P(\omega_k)$

We know $\sigma$

$p(x | \omega_i, \mu_i, ... \mu_k)$ = Prob density that an observation
from class $\omega_i$ would have value $x$
given class means $\mu_1... \mu_k$

Can we write an expression for that?

Yes: The standard multivariate Gaussian using mean $\mu_i$

Clustering with Gaussian Mixtures: Slide 16

# likelihoods in unsupervised case

We have $x_1, x_2 \ldots x_R$
We have $P(\omega_1), \ldots, P(\omega_k)$.  We have σ.
We can define, for any $x$, $p(x | \omega_i, \mu_1, \mu_2 .. \mu_k)$

Can we define $p(x | \mu_1, \mu_2 .. \mu_k)$ ?

Can we define $p(x_1, x_1, .. x_n | \mu_1, \mu_2 .. \mu_k)$ ?

---

# likelihoods in unsupervised case

We have $x_1, x_2 \ldots x_R$
We have $P(\omega_1), \ldots, P(\omega_k)$.  We have σ.
We can define, for any $x$, $p(x | \omega_i, \mu_1, \mu_2 .. \mu_k)$

Can we define $p(x | \mu_1, \mu_2 .. \mu_k)$ ?

Yes: A weighted sum of multivariate Gaussians,
where the weighting of the i[th] component is $P(\omega_i)$

Can we define $p(x_1, x_1, .. x_n | \mu_1, \mu_2 .. \mu_k)$ ?

Yes, if we assume the x's were drawn independently

# 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; $\omega_1$ and $\omega_2$

$P(\omega_1) = 1/3 \quad P(\omega_2) = 2/3 \quad \sigma = 1$ .

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

Clustering with Gaussian Mixtures: Slide 19

---

# Duda & Hart's Example



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

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

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

  * corresponds to switching $\omega_1$ and $\omega_2$.

Clustering with Gaussian Mixtures: Slide 20

**10**

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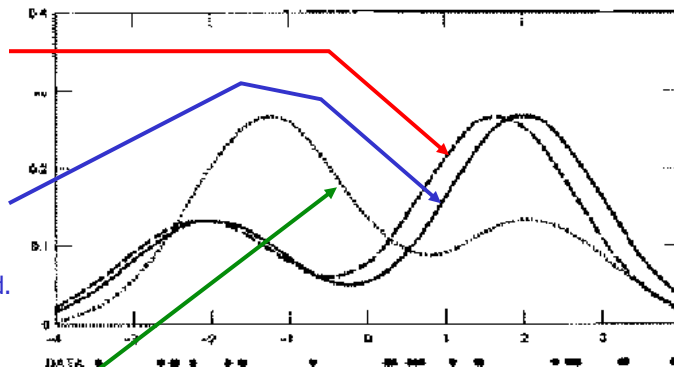- In this example unsupervised is almost as good as supervised. If the $x_1$ .. $x_{25}$ are given the class which was used to learn them, then the results are ($\mu_1$=-2.176, $\mu_2$=1.684). Unsupervised got ($\mu_1$=-2.13, $\mu_2$=1.668).

Clustering with Gaussian Mixtures: Slide 21

---

# Finding the max likelihood μ₁,μ₂..μₖ

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
    Slow but doable
- Use a much faster, cuter, and recently very popular method…

Clustering with Gaussian Mixtures: Slide 22

Expectation
Maximization

# 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 inference of Hidden Markov Models.

# Silly Example

Let events be "grades in a class"

| | | |
|---|---|---|
| $w_1$ = Gets an A | P(A) = ½ | |
| $w_2$ = Gets a  B | P(B) = μ | |
| $w_3$ = Gets a  C | P(C) = 2μ | |
| $w_4$ = Gets a  D | P(D) = ½-3μ | |

(Note  $0 \leq μ \leq 1/6$)

Assume we want to estimate μ 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 μ given a,b,c,d ?

Clustering with Gaussian Mixtures: Slide 25

---

# Trivial Statistics

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

$P(\,a,b,c,d\mid μ) = K(½)^a(μ)^b(2μ)^c(½\text{-}3μ)^d$

$\log P(\,a,b,c,d\mid μ) = \log K + a\log ½ + b\log μ + c\log 2μ + d\log (½\text{-}3μ)$

FOR  MAX  LIKE  μ, SET  $\dfrac{\partial \log P}{\partial μ} = 0$

$\dfrac{\partial \log P}{\partial μ} = \dfrac{b}{μ} + \dfrac{2c}{2μ} - \dfrac{3d}{1/2 - 3μ} = 0$

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

So if  class  got

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

Max  like  $μ = \dfrac{1}{10}$

*Boring, but true!*

Clustering with Gaussian Mixtures: Slide 26

# 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 $\mu$ now?

REMEMBER

$P(A) = \frac{1}{2}$

$P(B) = \mu$

$P(C) = 2\mu$

$P(D) = \frac{1}{2}\text{-}3\mu$

---

# 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 $\mu$ now?

REMEMBER

$P(A) = \frac{1}{2}$

$P(B) = \mu$

$P(C) = 2\mu$

$P(D) = \frac{1}{2}\text{-}3\mu$

We can answer this question circularly:

**EXPECTATION**

If we know the value of $\mu$ we could compute the expected value of $a$ and $b$

Since the ratio a:b should be the same as the ratio ½ : $\mu$

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

**MAXIMIZATION**

If we know the true values of $a$ and $b$ we could compute the maximum likelihood value of $\mu$

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

# E.M. for our Trivial Problem

We begin with a guess for μ

We iterate between EXPECTATION and MAXIMIZATION to improve our estimates of μ and *a* and *b*.

Define    μ(t)  the estimate of μ on the t'th iteration

          b(t)  the estimate of *b* on t'th iteration

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

$$b(t) = \frac{\mu(t)h}{\frac{1}{2} + \mu(t)} = \mathrm{E}[b \mid \mu(t)]$$

**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.**

     Clustering with Gaussian Mixtures: Slide 29

---

# 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]
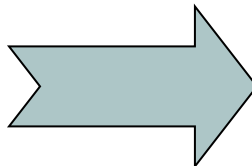
So it must therefore converge   [OBVIOUS]

In our example, suppose we had

   h = 20
   c = 10
   d = 10
   μ(0) = 0

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

| t | μ(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 |

     Clustering with Gaussian Mixtures: Slide 30

# Back to Unsupervised Learning of Gaussian Mixture Models

Remember:

   We have unlabeled data $x_1\ x_2\ ...\ x_R$

   We know there are k classes

   We know $P(\omega_1)\ P(\omega_2)\ P(\omega_3)\ ...\ P(\omega_k)$

   We <u>don't</u> know $\boldsymbol{\mu}_1\ \boldsymbol{\mu}_2\ ..\ \boldsymbol{\mu}_k$

We can write p( data | $\boldsymbol{\mu}_1 .... \boldsymbol{\mu}_k$)

$$= p(\mathbf{x}_1 ... \mathbf{x}_R | \mu_1 ... \mu_k)$$

$$= \prod_{i=1}^{R} p(\mathbf{x}_i | \mu_1 ... \mu_k)$$

$$= \prod_{i=1}^{R} \sum_{j=1}^{k} p(\mathbf{x}_i | \omega_j, \mu_1 ... \mu_k) P(\omega_j)$$

$$= \prod_{i=1}^{R} \sum_{j=1}^{k} K \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mu_j\|^2\right) P(\omega_j)$$

                Clustering with Gaussian Mixtures: Slide 31

---

# E.M. for GMMs

For Max likelihood we know $\quad \dfrac{\partial}{\partial \mu_i} \log \operatorname{Pr}ob(\text{data}|\mu_1 ... \mu_k) = 0$

Some wild'n'crazy algebra turns this into :"For Max likelihood, for each j,

$$\mu_j = \frac{\displaystyle\sum_{i=1}^{R} P(\omega_j | \mathbf{x}_i, \mu_1 ... \mu_k)\, \mathbf{x}_i}{\displaystyle\sum_{i=1}^{R} P(\omega_j | \mathbf{x}_i, \mu_1 ... \mu_k)}$$

This is  n  nonlinear equations in $\boldsymbol{\mu}_j$'s."

If, for each $\mathbf{x}_i$ we knew that for each $\omega_j$ the prob that $\mathbf{x}_i$ was in class $\omega_j$ is $P(\omega_j | x_i, \mu_1 ... \mu_k)$ ...  then we could easily compute $\boldsymbol{\mu}_j$.

If we knew each $\mu_j$ then we could easily compute $P(\omega_j | \mathbf{x}_i, \mu_1 ... \mu_j)$ for each $\omega_j$ and $\mathbf{x}_i$.

                          ...I feel an EM experience coming on!!

                Clustering with Gaussian Mixtures: Slide 32

# E.M. for GMMs

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

$$\lambda_t = \{\, \mu_1(t),\, \mu_2(t) \ldots \mu_c(t)\,\}$$

E-step

Compute "expected" classes of all datapoints for each class

$$P\!\left(\omega_i \middle| \mathbf{x}_k, \lambda_t\right) = \frac{p\!\left(\mathbf{x}_k \middle| \omega_i, \lambda_t\right) P\!\left(\omega_i \middle| \lambda_t\right)}{p\!\left(\mathbf{x}_k \middle| \lambda_t\right)} = \frac{p\!\left(\mathbf{x}_k \middle| \omega_i, \mathbf{\mu}_i(t), \sigma^2 \mathbf{I}\right) p_i(t)}{\sum_{j=1}^{c} p\!\left(\mathbf{x}_k \middle| \omega_j, \mathbf{\mu}_j(t), \sigma^2 \mathbf{I}\right) p_j(t)}$$

*Just evaluate a Gaussian at $x_k$*

M-step.

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

$$\mu_i(t+1) = \frac{\sum_k P\!\left(\omega_i \middle| \mathbf{x}_k, \lambda_t\right) \mathbf{x}_k}{\sum_k P\!\left(\omega_i \middle| \mathbf{x}_k, \lambda_t\right)}$$

---

# E.M. Convergence



- As with all EM procedures, 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

# E.M. for General GMMs

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

$$\lambda_t = \{\, \mu_1(t),\ \mu_2(t)\ \dots\ \mu_c(t),\ \Sigma_1(t),\ \Sigma_2(t)\ \dots\ \Sigma_c(t),\ p_1(t),\ p_2(t)\ \dots\ p_c(t)\,\}$$

> $p_i(t)$ is shorthand for estimate of $P(\omega_i)$ on t'th iteration

**E-step**

Compute "expected" classes of all datapoints for each class

> Just evaluate a Gaussian at $x_k$

$$P(\omega_i | x_k, \lambda_t) = \frac{p(\mathbf{x}_k | \omega_i, \lambda_t) P(\omega_i | \lambda_t)}{p(\mathbf{x}_k | \lambda_t)} = \frac{p(\mathbf{x}_k | \omega_i, \mathbf{\mu}_i(t), \mathbf{\Sigma}_i(t)) p_i(t)}{\sum_{j=1}^{c} p(\mathbf{x}_k | \omega_j, \mathbf{\mu}_j(t), \mathbf{\Sigma}_j(t)) p_j(t)}$$

**M-step.**

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

$$\mu_i(t+1) = \frac{\sum_k P(\omega_i | \mathbf{x}_k, \lambda_t) \mathbf{x}_k}{\sum_k P(\omega_i | \mathbf{x}_k, \lambda_t)} \qquad \mathbf{\Sigma}_i(t+1) = \frac{\sum_k P(\omega_i | \mathbf{x}_k, \lambda_t)[\mathbf{x}_k - \mathbf{\mu}_i(t+1)][\mathbf{x}_k - \mathbf{\mu}_i(t+1)]^T}{\sum_k P(\omega_i | \mathbf{x}_k, \lambda_t)}$$

$$p_i(t+1) = \frac{\sum_k P(\omega_i | \mathbf{x}_k, \lambda_t)}{R}$$

> $R$ = #records
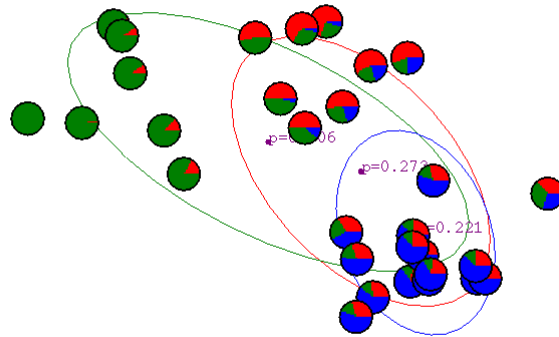
Clustering with Gaussian Mixtures: Slide 35

---

# Gaussian Mixture Example: Start



*Advance apologies: in Black and White this example will be incomprehensible*

Clustering with Gaussian Mixtures: Slide 36

After first
iteration

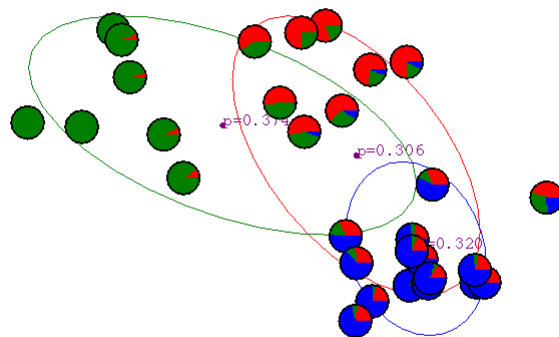Clustering with Gaussian Mixtures: Slide 37



After 2nd
iteration

Clustering with Gaussian Mixtures: Slide 38

**19**

After 3rd iteration

p=0.34
p=0.307

Copyright © 2001, Andrew W. Moore

Clustering with Gaussian Mixtures: Slide 39



After 4th iteration

p=0.331
p=0.288

Copyright © 2001, Andrew W. Moore
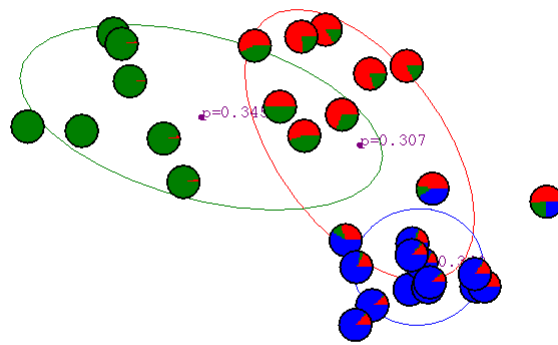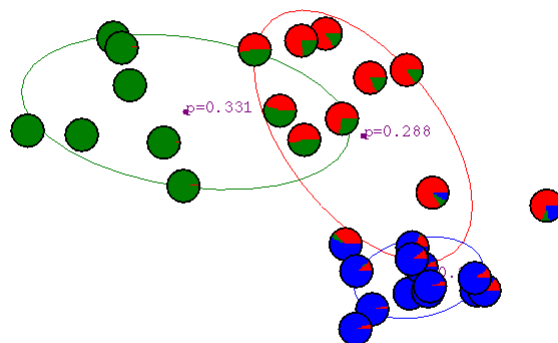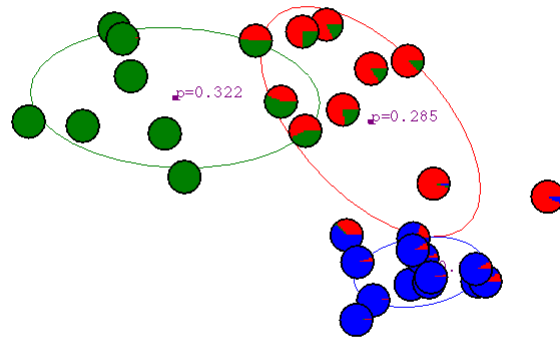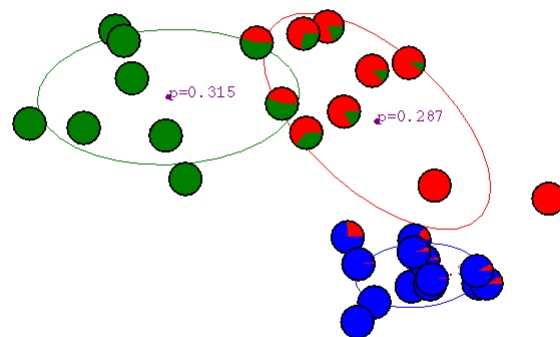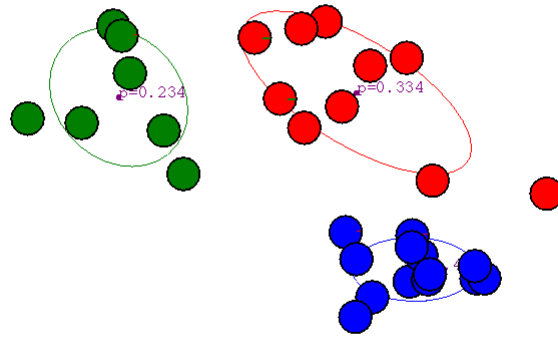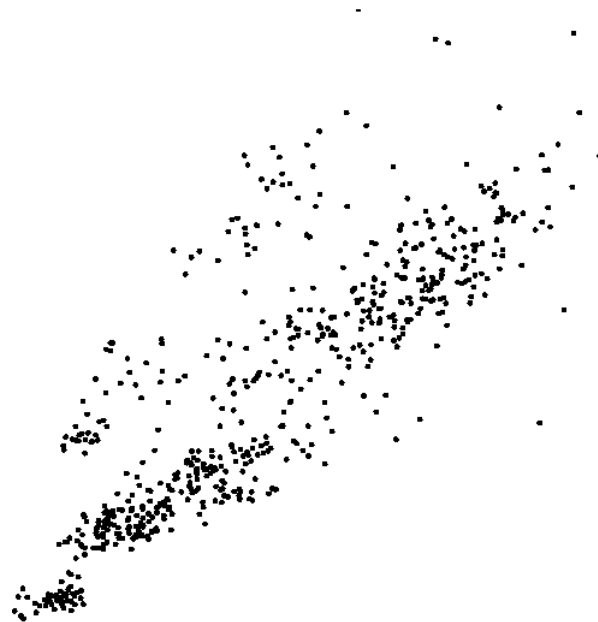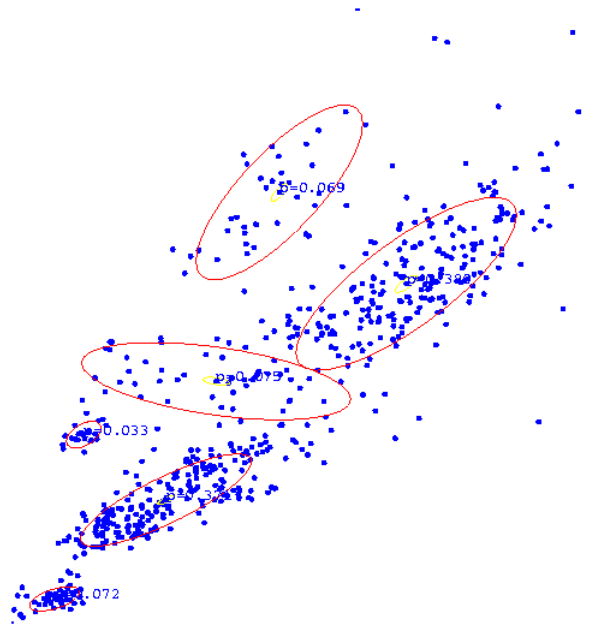
Clustering with Gaussian Mixtures: Slide 40

After 5th iteration

p=0.322   p=0.285

After 6th iteration

p=0.315   p=0.287

After 20th iteration

Clustering with Gaussian Mixtures: Slide 43



Some Bio Assay data

Clustering with Gaussian Mixtures: Slide 44

# GMM clustering of the assay data
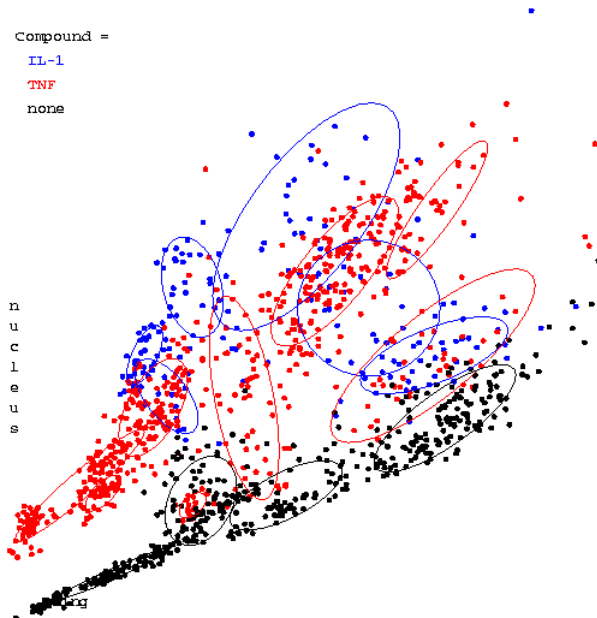
# Resulting Density Estimator

Three classes of assay
(each learned with it's own mixture model)
(Sorry, this will again be semi-useless in black and white)
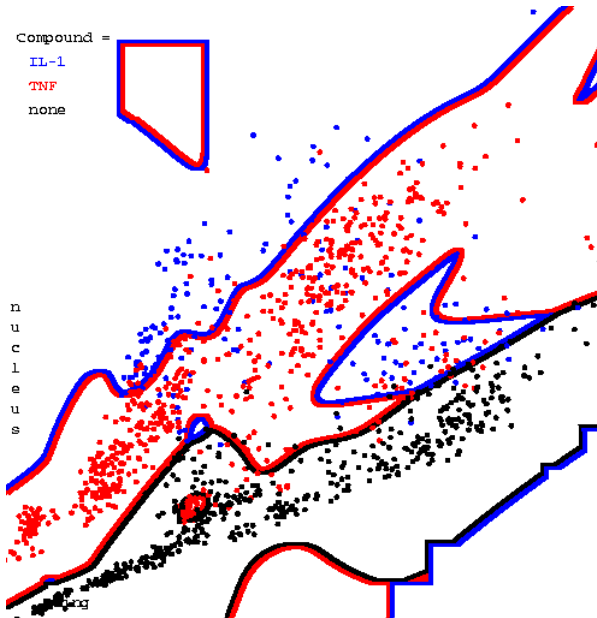
Compound =
IL-1
TNF
none

nucleus

Copyright © 2001, Andrew W. Moore

Clustering with Gaussian Mixtures: Slide 47



Resulting Bayes Classifier
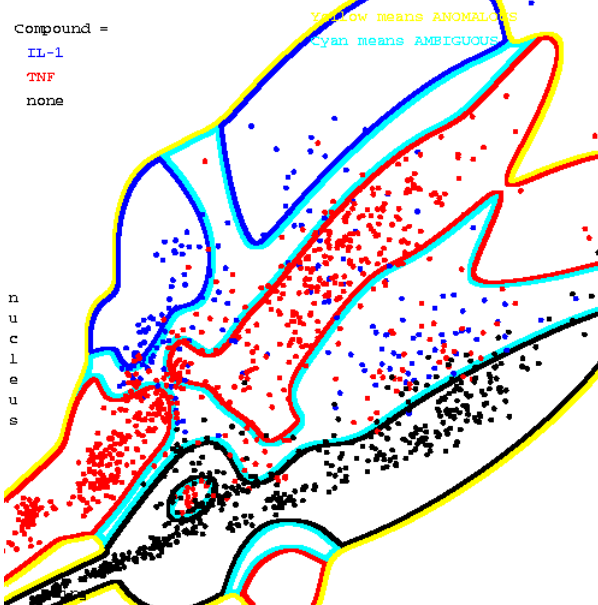
Compound =
IL-1
TNF
none

nucleus

Copyright © 2001, Andrew W. Moore

Clustering with Gaussian Mixtures: Slide 48

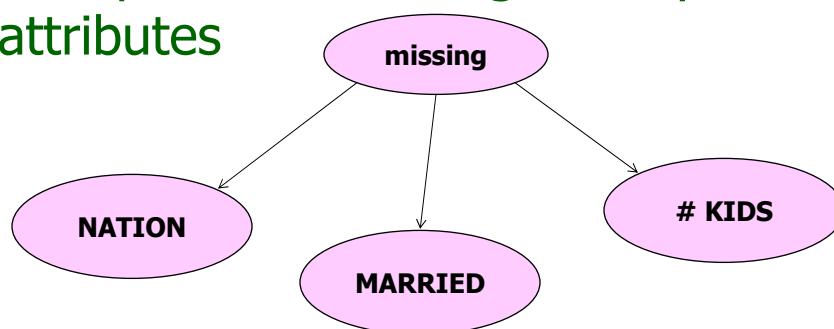Resulting Bayes Classifier, using posterior probabilities to alert about ambiguity and anomalousness

Compound =
IL-1
TNF
none

Yellow means ANOMALOUS
Cyan means AMBIGUOUS

n u c l e u s

**Yellow means anomalous**

**Cyan means ambiguous**

Clustering with Gaussian Mixtures: Slide 49

---

# Unsupervised learning with symbolic attributes



missing

NATION

MARRIED

# KIDS

It's just a "learning Bayes net with known structure but hidden values" problem.

Can use Gradient Descent.

EASY, fun exercise to do an EM formulation for this case too.

Clustering with Gaussian Mixtures: Slide 50

# Final Comments

- Remember, E.M. can get stuck in local minima, and empirically it <u>DOES</u>.

- Our unsupervised learning example assumed $P(\omega_i)$'s known. Easy to relax this.

- It's possible to do Bayesian unsupervised learning instead of max. likelihood.

- There are other algorithms for unsupervised learning. We'll visit K-means soon. Hierarchical clustering is also interesting.

- Neural-net algorithms called "competitive learning" turn out to have interesting parallels with the EM method we saw.

# What you should know

- How to "learn" maximum likelihood parameters (locally max. like.) in the case of unlabeled data.

- Be happy with this kind of probabilistic analysis.

- Understand the two examples of E.M. given in these notes.

# Other unsupervised learning methods

- K-means (see next lecture)
- Hierarchical clustering (e.g. Minimum spanning trees)
- Principal Component Analysis
    simple, useful tool

- Non-linear PCA
    Neural Auto-Associators
    Locally weighted PCA
    Others…

Clustering with Gaussian Mixtures: Slide 53