# Clustering with Gaussian Mixtures

**Andrew W. Moore**

**Associate Professor**

**School of Computer Science**

**Carnegie Mellon University**

www.cs.cmu.edu/~awm
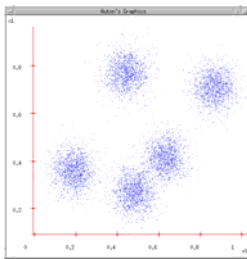
awm@cs.cmu.edu

412-268-7599

　　　　　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 all the $P(w_j)$'s ."
- So far, looks straightforward.
  "I need a maximum likelihood estimate of the $\mu_i$'s ."
- No problem:
  "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

---

## Some data from a GMM



　　　Clustering with Gaussian Mixtures: Slide 3

---

## The GMM assumption

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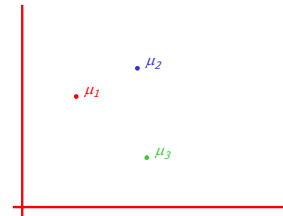


　　　Clustering with Gaussian Mixtures: Slide 4

---

## The 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 $\sigma^2 I$

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



　　　Clustering with Gaussian Mixtures: Slide 5

---

## The GMM assumption

- There are k components. The i'th component is called $\omega_i$
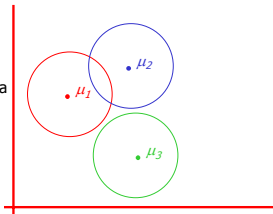- Component $\omega_i$ has an associated mean vector $\mu_i$
- Each component generates data from a Gaussian with mean $\mu_i$ and covariance matrix $\sigma^2 I$

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

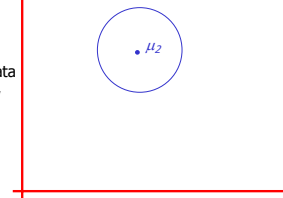1. Pick a component at random. Choose component i with probability $P(\omega_i)$.
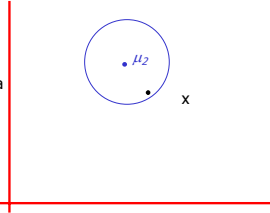


　　　Clustering with Gaussian Mixtures: Slide 6

## The 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 $\sigma^2 I$

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 ~ $N(\mu_i, \sigma^2 I)$

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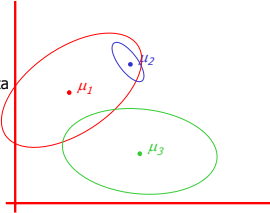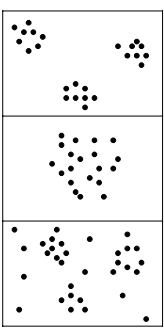2. Datapoint ~ $N(\mu_i, \Sigma_i)$

## 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 $x_1, x_2, \ldots x_N$
We know $P(w_1) P(w_2) \ldots P(w_k)$
We know $\sigma$

$P(x|w_i, \mu_i, \ldots \mu_k)$ = Prob that an observation from class $w_i$ would have value $x$ given class means $\mu_1 \ldots \mu_x$

Can we write an expression for that?

## likelihoods in unsupervised case

We have $x_1 x_2 \ldots x_n$
We have $P(w_1) \ldots P(w_k)$. We have $\sigma$.
We can define, for any $x$, $P(x|w_i, \mu_1, \mu_2 \ldots \mu_k)$

Can we define $P(x \mid \mu_1, \mu_2 \ldots \mu_k)$ ?

Can we define $P(x_1, x_2, \ldots x_n \mid \mu_1, \mu_2 \ldots \mu_k)$ ?

[YES, IF WE ASSUME THE $X_i$'S WERE DRAWN INDEPENDENTLY]

## Unsupervised Learning: Mediumly Good News

We now have a procedure s.t. if you give me a guess at $\mu_1, \mu_2 \ldots \mu_k$,
I can tell you the prob of the unlabeled data given those $\mu$'s.

Suppose $x$'s are 1-dimensional.

There are two classes; $w_1$ and $w_2$
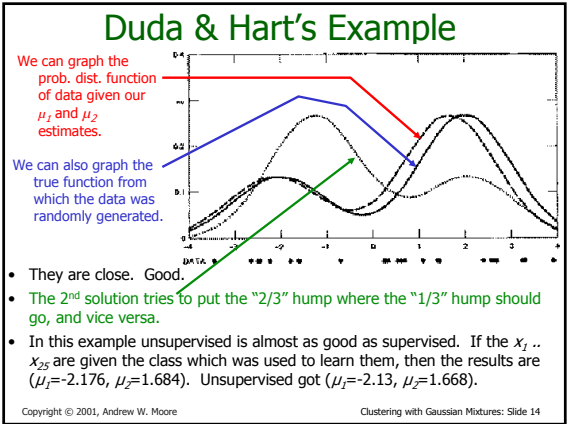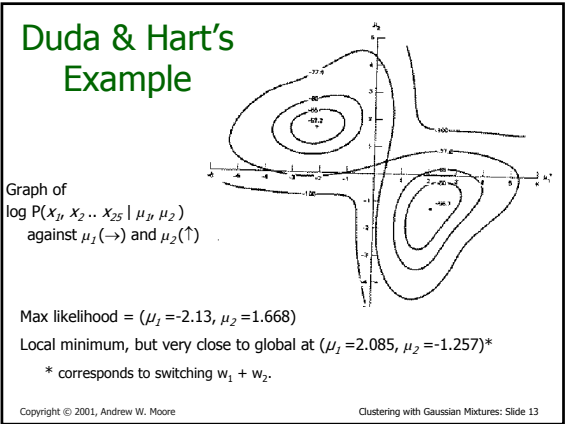
$P(w_1) = 1/3$    $P(w_2) = 2/3$    $\sigma = 1$.

There are 25 unlabeled datapoints

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

**(From Duda and Hart)**

DATA SCATTERGRAM

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

  * corresponds to switching $w_1 + w_2$.

Clustering with Gaussian Mixtures: Slide 13

---

## 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).
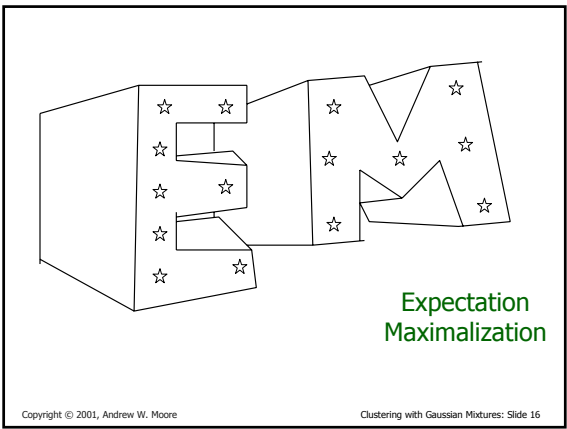
Clustering with Gaussian Mixtures: Slide 14

---

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

We can compute  P( data $\mid \mu_1 \mu_2 .. \mu_k$)
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 15

---



Expectation Maximalization

Clustering with Gaussian Mixtures: Slide 16

---

## 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 (future lecture).

Clustering with Gaussian Mixtures: Slide 17

---

## 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  $0 \leq \mu \leq 1/6$)
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 ?

Clustering with Gaussian Mixtures: Slide 18

## Computing

$P(A) = \tfrac{1}{2}$ $\quad P(B) = \mu$ $\quad P(C) = 2\mu$ $\quad P(D) = \tfrac{1}{2}-3\mu$

$P(a,b,c,d \mid \mu) = K(\tfrac{1}{2})^a(\mu)^b(2\mu)^c(\tfrac{1}{2}-3\mu)^d$

$\log P(a,b,c,d \mid \mu) = \log K + a\log \tfrac{1}{2} + b\log \mu + c\log 2\mu + d\log(\tfrac{1}{2}-3\mu)$

FOR MAX LIKE $\mu$, SET $\dfrac{\partial \mathrm{Log}P}{\partial \mu} = 0$

$\dfrac{\partial \mathrm{Log}P}{\partial \mu} = \dfrac{b}{\mu} + \dfrac{2c}{2\mu} - \dfrac{3d}{1/2 - 3\mu} = 0$

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

Clustering with Gaussian Mixtures: Slide 19

---

## 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
$P(A) = \tfrac{1}{2}$
$P(B) = \mu$
$P(C) = 2\mu$
$P(D) = \tfrac{1}{2}-3\mu$

Clustering with Gaussian Mixtures: Slide 20

---

## 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?

We can answer this question circularly:

REMEMBER
$P(A) = \tfrac{1}{2}$
$P(B) = \mu$
$P(C) = 2\mu$
$P(D) = \tfrac{1}{2}-3\mu$

**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 be the same as the ratio ½ : µ

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

**MAXIMIZATION**

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

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

Clustering with Gaussian Mixtures: Slide 21

---

## E.M. for our Trivial Problem

REMEMBER
$P(A) = \tfrac{1}{2}$
$P(B) = \mu$
$P(C) = 2\mu$
$P(D) = \tfrac{1}{2}-3\mu$

We begin with a guess for $\mu$

We iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of $\mu$ and $a$ and $b$.

Define $\mu(t)$ the estimate of $\mu$ on the t'th iteration
$b(t)$ the estimate of $b$ on t'th iteration

$\mu(0) = $ initial guess

$b(t) = \dfrac{\mu(t)h}{\tfrac{1}{2}+\mu(t)} = E[b \mid \mu(t)]$ — **E-step**

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

$= $ max like est of $\mu$ given $b(t)$
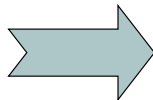
Continue iterating until converged.
**Good news:** Converging to local optimum is assured.
**Bad news:** I said "local" optimum.

Clustering with Gaussian Mixtures: Slide 22

---

## E.M. Convergence

- Convergence proof based on fact that Prob(data | $\mu$) 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
$h = 20$
$c = 10$
$d = 10$
$\mu(0) = 0$

| 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 |

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

Clustering with Gaussian Mixtures: Slide 23

---

## Back to Unsupervised Learning of GMMs

Remember:
We have unlabeled data $x_1\, x_2 \ldots x_R$
We know there are k classes
We know $P(w_1)\, P(w_2)\, P(w_3) \ldots P(w_k)$
We <u>don't</u> know $\mu_1\, \mu_2 \ldots \mu_k$

We can write $P(\text{data} \mid \mu_1 \ldots \mu_k)$

$= p(x_1 \ldots x_R \mid \mu_1 \ldots \mu_k)$

$= \prod_{i=1}^{R} p(x_i \mid \mu_1 \ldots \mu_k)$

$= \prod_{i=1}^{R} \sum_{j=1}^{k} p(x_i \mid w_j, \mu_1 \ldots \mu_k) P(w_j)$

$= \prod_{i=1}^{R} \sum_{j=1}^{k} K \exp\left(-\dfrac{1}{2\sigma^2}(x_i - \mu_j)^2\right) P(w_j)$

Clustering with Gaussian Mixtures: Slide 24

**4**

## E.M. for GMMs

For Max likelihood we know $\dfrac{\partial}{\partial \mu_i}\log\mathrm{Pr\,ob}\big(\mathrm{data}\,|\,\mu_1\ldots\mu_k\big)=0$

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

$$\mu_j=\frac{\sum_{i=1}^{R}P\big(w_j\,\big|\,x_i,\mu_1\ldots\mu_k\big)x_i}{\sum_{i=1}^{R}P\big(w_j\,\big|\,x_i,\mu_1\ldots\mu_k\big)}$$

This is n  nonlinear equations in $\mu_j$'s."

If, for each $\mathbf{x}_i$ we knew that for each $w_j$ the prob that $\mu_j$ was in class $w_j$ is $P(w_j|x_i,\mu_1\ldots\mu_k)$   Then… we would easily compute $\mu_j$.

If we knew each $\mu_j$ then we could easily compute $P(w_j|x_i,\mu_1\ldots\mu_j)$ for each $w_j$ and $x_i$.

…I feel an EM experience coming on!!

Clustering with Gaussian Mixtures: Slide 25

---

## 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

*Just evaluate a Gaussian at $x_k$*

$$P\big(w_i\,|\,x_k,\lambda_t\big)=\frac{p\big(x_k\,|\,w_i,\lambda_t\big)P\big(w_i\,|\,\lambda_t\big)}{p\big(x_k\,|\,\lambda_t\big)}=\frac{p\big(x_k\,|\,w_i,\mu_i(t),\sigma^2\mathbf{I}\big)p_i(t)}{\sum_{j=1}^{c}p\big(x_k\,|\,w_j,\mu_j(t),\sigma^2\mathbf{I}\big)p_j(t)}$$

**M-step.**

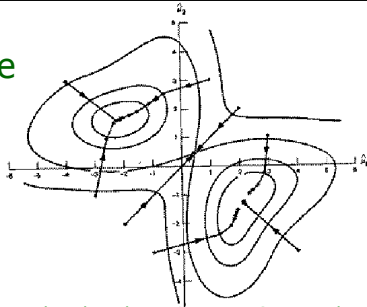Compute Max. like $\mu$ given our data's class membership distributions

$$\mu_i(t+1)=\frac{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)x_k}{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)}$$

Clustering with Gaussian Mixtures: Slide 26

---

## 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

Clustering with Gaussian Mixtures: Slide 27

---

## E.M. for General GMMs

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

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

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

**E-step**

Compute "expected" classes of all datapoints for each class

*Just evaluate a Gaussian at $x_k$*

$$P\big(w_i\,|\,x_k,\lambda_t\big)=\frac{p\big(x_k\,|\,w_i,\lambda_t\big)P\big(w_i\,|\,\lambda_t\big)}{p\big(x_k\,|\,\lambda_t\big)}=\frac{p\big(x_k\,|\,w_i,\mu_i(t),\Sigma_i(t)\big)p_i(t)}{\sum_{j=1}^{c}p\big(x_k\,|\,w_j,\mu_j(t),\Sigma_j(t)\big)p_j(t)}$$

**M-step.**

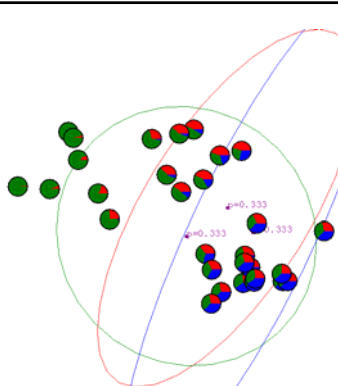Compute Max. like $\mu$ given our data's class membership distributions

$$\mu_i(t+1)=\frac{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)x_k}{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)}\qquad \Sigma_i(t+1)=\frac{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)\big[x_k-\mu_i(t+1)\big]\big[x_k-\mu_i(t+1)\big]^T}{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)}$$

$$p_i(t+1)=\frac{\sum_k P\big(w_i\,|\,x_k,\lambda_t\big)}{R}$$

*R = #records*

Clustering with Gaussian Mixtures: Slide 28

---

## Gaussian Mixture Example: Start



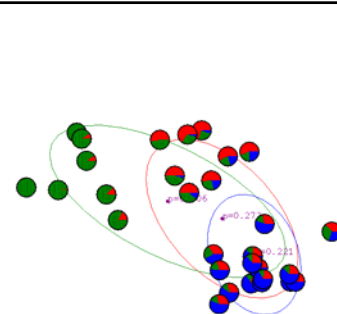*Advance apologies: in Black and White this example will be incomprehensible*

Clustering with Gaussian Mixtures: Slide 29

---

## After first iteration

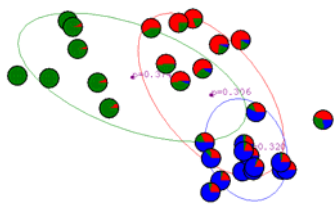

Clustering with Gaussian Mixtures: Slide 30

After 2nd iteration

Clustering with Gaussian Mixtures: Slide 31



After 3rd iteration

Clustering with Gaussian Mixtures: Slide 32



After 4th iteration

Clustering with Gaussian Mixtures: Slide 33



After 5th iteration

Clustering with Gaussian Mixtures: Slide 34



After 6th iteration

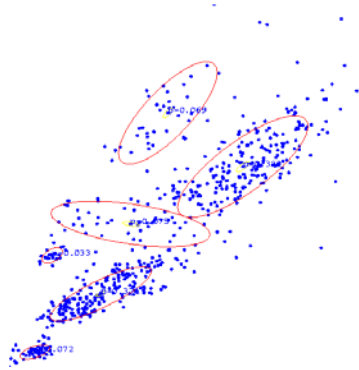Clustering with Gaussian Mixtures: Slide 35



After 20th iteration

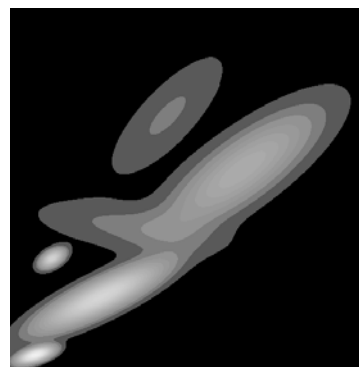Clustering with Gaussian Mixtures: Slide 36

Some Bio Assay data

Copyright © 2001, Andrew W. Moore — Clustering with Gaussian Mixtures: Slide 37


GMM clustering of the assay data

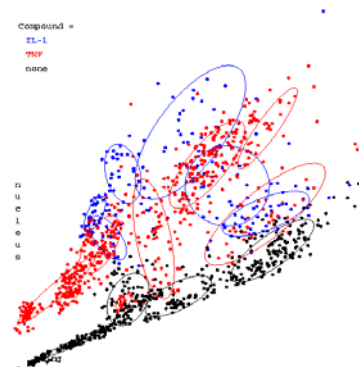Copyright © 2001, Andrew W. Moore — Clustering with Gaussian Mixtures: Slide 38


Resulting Density Estimator

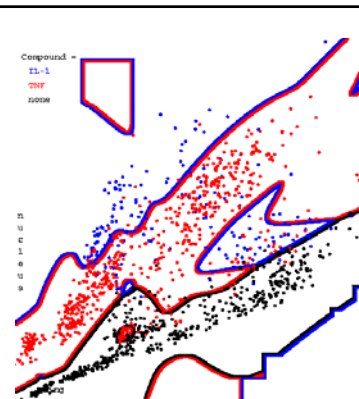Copyright © 2001, Andrew W. Moore — Clustering with Gaussian Mixtures: Slide 39


Three classes of assay
(each learned with it's own mixture model)
(Sorry, this will again be semi-useless in black and white)

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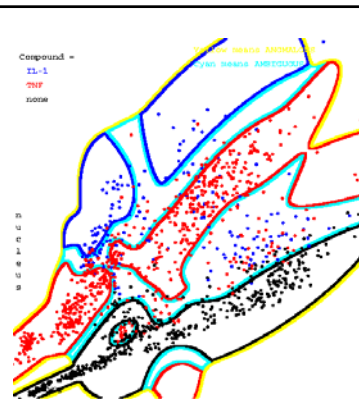

Resulting Bayes Classifier

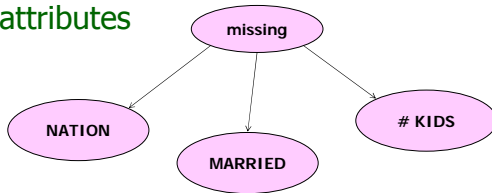Copyright © 2001, Andrew W. Moore — Clustering with Gaussian Mixtures: Slide 41


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

Yellow means anomalous

Cyan means ambiguous

Copyright © 2001, Andrew W. Moore — Clustering with Gaussian Mixtures: Slide 42

## Unsupervised learning with symbolic attributes



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 43

---

## Final Comments

- Remember, E.M. can get stuck in local minima, and empirically it DOES.
- Our unsupervised learning example assumed $P(w_i)$'s known, and variances fixed and 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.

Clustering with Gaussian Mixtures: Slide 44

---

## 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.

  For more info, see Duda + Hart. It's a great book. There's much more in the book than in your handout.

Clustering with Gaussian Mixtures: Slide 45

---

## Other unsupervised learning methods

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

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

Clustering with Gaussian Mixtures: Slide 46