# Unsupervised learning or Clustering –
# K-means
# Gaussian mixture models

*E.M.*

Machine Learning – 10701/15781
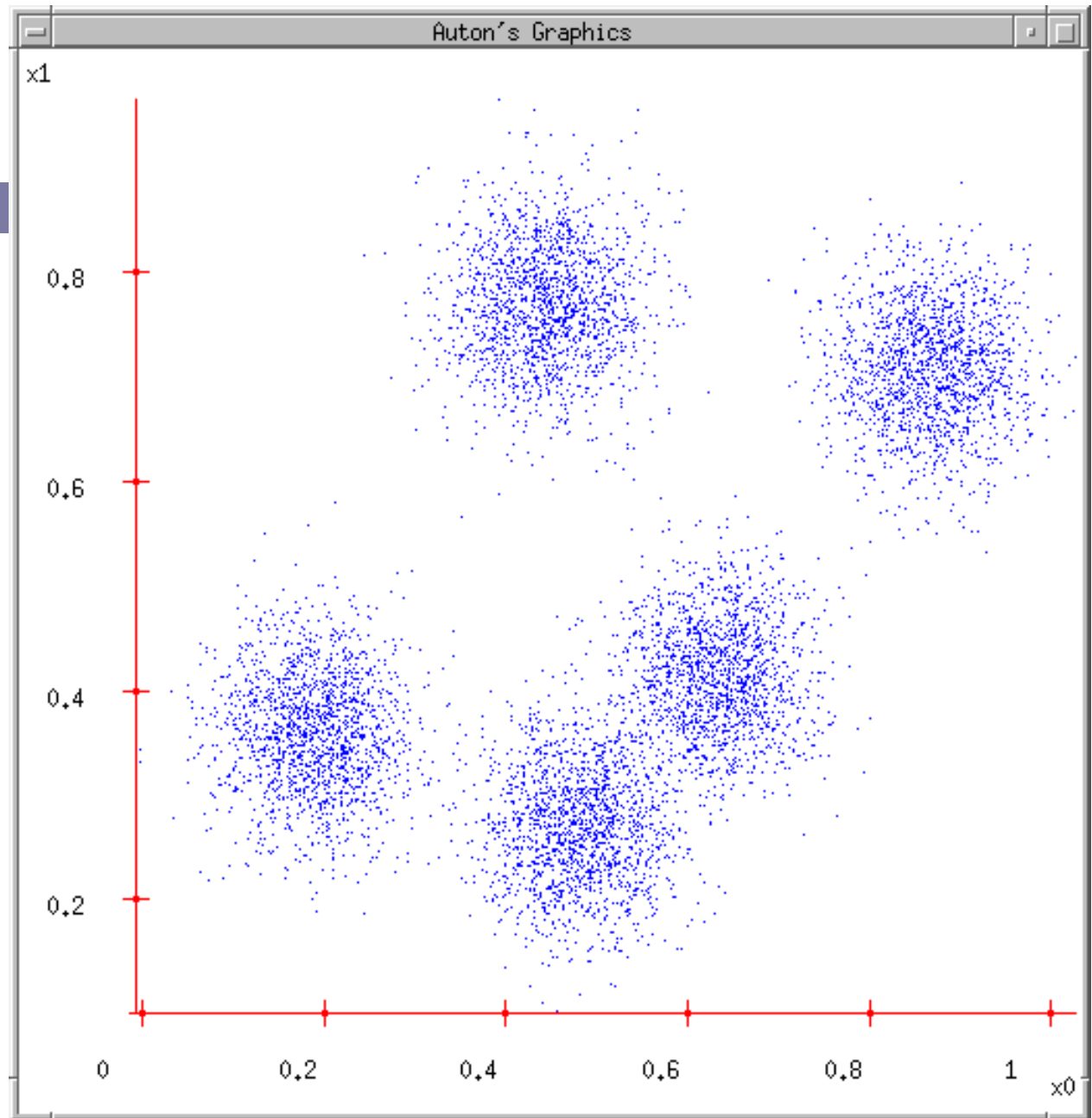
Carlos Guestrin

Carnegie Mellon University
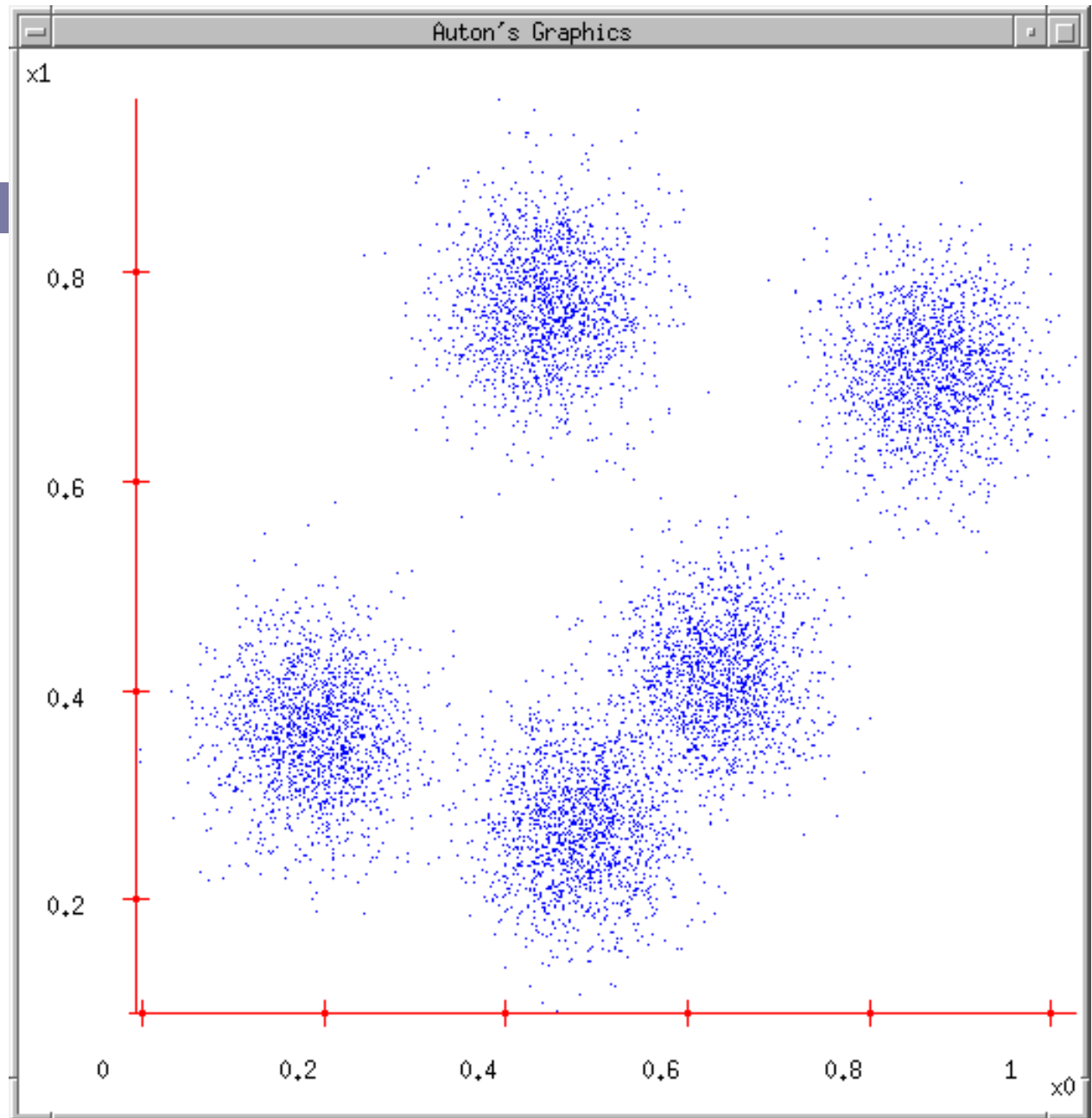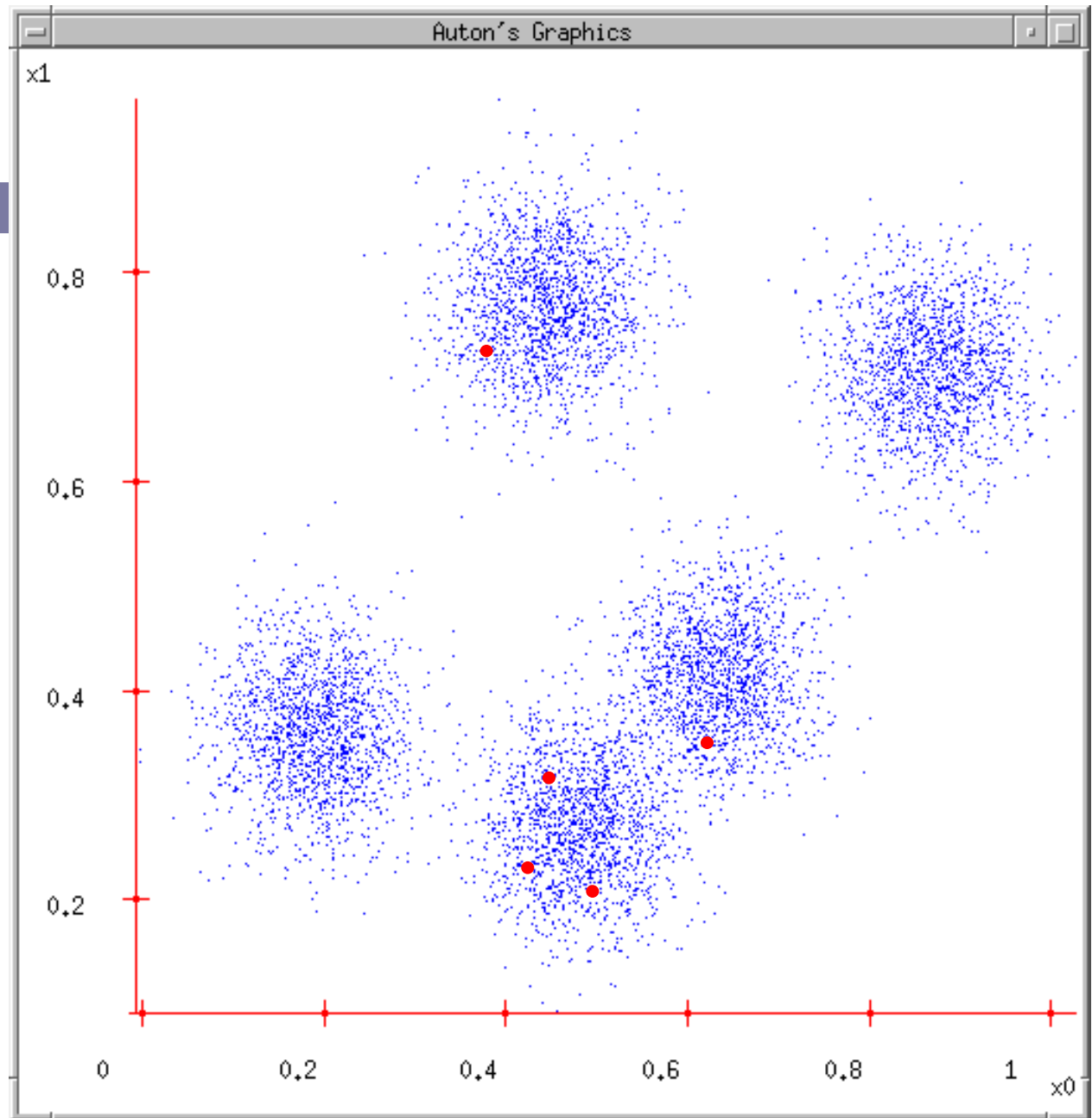
April 4th, 2007

# Some Data

# K-means

1. Ask user how many clusters they'd like.
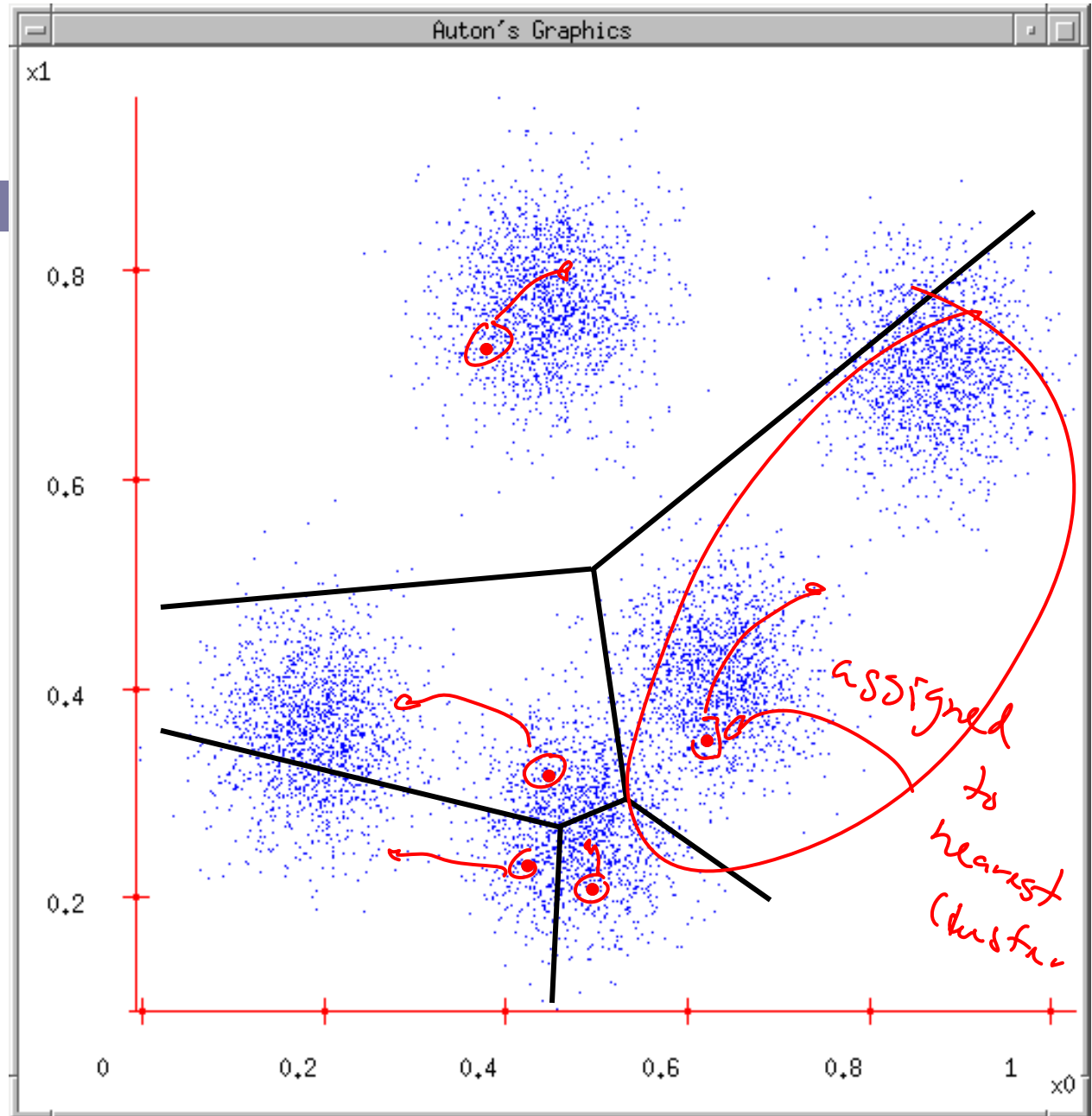   *(e.g. k=5)*

# K-means

1. Ask user how many clusters they'd like. *(e.g. k=5)*

2. Randomly guess k cluster Center locations

# 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. (Thus each Center "owns" a set of datapoints)



assigned to nearest cluster

# 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



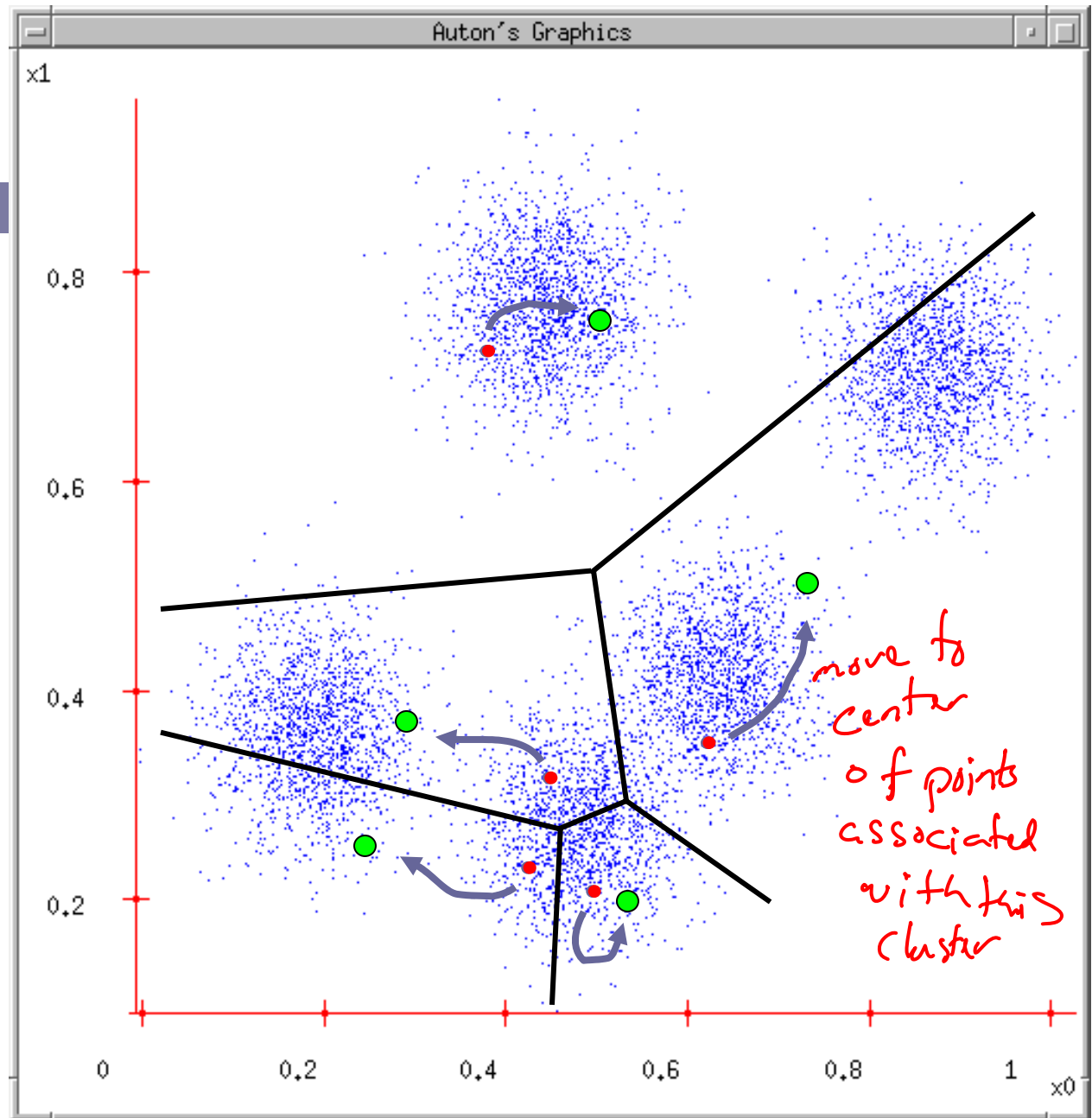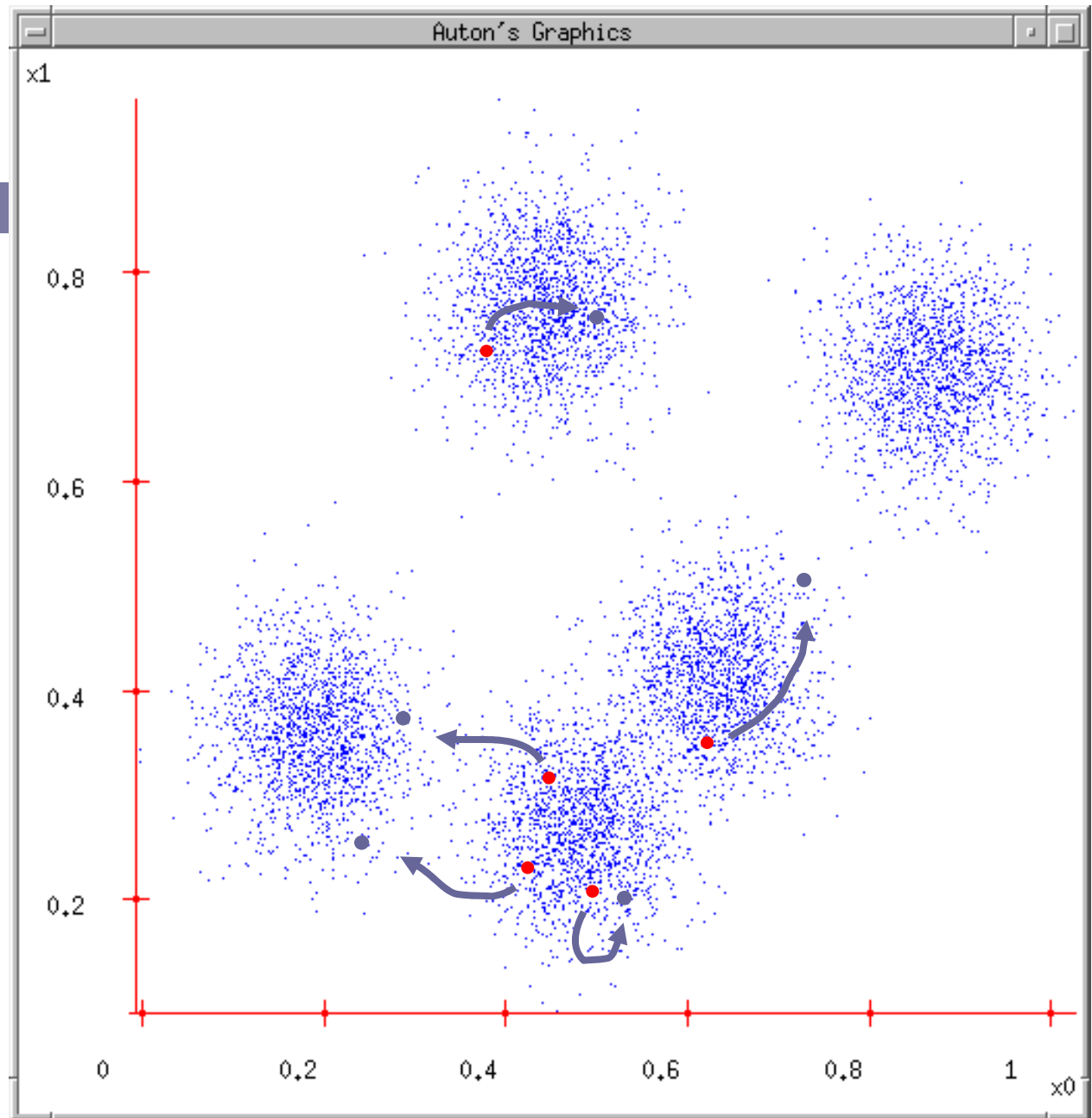move to center of points associated with this cluster

# 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.
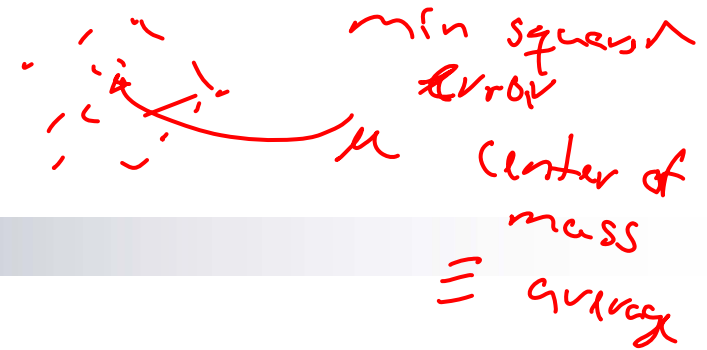
4.  Each Center finds the centroid of the points it owns...

5.  ...and jumps there

6.  ...Repeat until terminated!

©2005-2007 Carlos Guestrin

# K-means

*min squared error*
*center of mass*
*≡ average*

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

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

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

# What is K-means optimizing?

- Potential function F(μ,C) of centers μ and point allocations C:

  - $$F(\mu, C) = \sum_{j=1}^{m} ||\mu_{C(j)} - x_j||^2$$

  *distance between $x_j$ & $\mu_{C(j)}$ ← nearest cluster center of $x_j$*

- Optimal K-means:
  - $\min_\mu \min_C F(\mu, C)$

    *Center locations    assignment points to clusters*

# Does K-means converge??? Part 1

- Optimize potential function:

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

*clusters*

*sum over points in Cluster i*

- Fix μ, optimize C

set μ to μ̂ , opt. over C

$$\min_{C} \sum_{i=1}^{k} \sum_{C(j)=i} ||\hat{\mu}_i - x_j||^2 = \min_{C} \sum_{j=1}^{m} ||\hat{\mu}_{C(j)} - x_j||^2$$

$$= \sum_{j=1}^{m} \min_{C(j)} ||\hat{\mu}_{C(j)} - x_j||^2 \longleftarrow$$

because assign of C(j) is *choose independently* of C(ℓ)

independently pick

$$C(j) = \arg\min_{i} ||\hat{\mu}_i - x_j||^2$$

~~Reclassify~~ Classification

# Does K-means converge??? Part 2

- Optimize potential function:

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

- Fix C, optimize μ $\quad$ set $C$ to $\hat{C}$

$$\min_{\mu} \sum_{i=1}^{k} \sum_{j:\hat{C}(j)=i} ||\mu_i - x_j||^2 = \quad \text{can pick centers } \mu_i \text{ independently}$$

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

recenter $\qquad \mu_i = \text{argmin}_{\mu_i} \sum_{j:\hat{C}(j)=i} ||\mu_i - x_j||^2$

# Coordinate descent algorithms

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

$\geq 0$

- Want: $\min_a \min_b F(a,b)$
- Coordinate descent:
  - fix a, minimize b
  - fix b, minimize a
  - repeat

fix

- Converges!!!

$\exists c \; \forall_{a,b} \; F(a,b) \geq c$  fix $b$

  - if F is bounded
  - to a (often good) local optimum
    - as we saw in applet (play with it!)

$b$

fix $\mu$   OPT C

fix C   OPT $\mu$

- K-means is a coordinate descent algorithm!

# (One) bad case for k-means

- Clusters may overlap

- Some clusters may be "wider" than others



definitely $C_1$

0/1
hard decision of assignment points to clusters

mixture

definitely $C_2$

K-means can't do this

# Gaussian Bayes Classifier Reminder

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

class mean

class covariance

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

prior

Gaussian likelihood

# Predicting wealth from age



wealth = poor
(prior = 0.760718)

| 1 | mean | cov |
|---|------|-----|
| age | 37.374 | 198.935 |

wealth = rich
(prior = 0.239282)

| 1 | mean | cov |
|---|------|-----|
| age | 44.7727 | 111.618 |

$P(X | y = poor)$

$P(X | y = rich)$

# Predicting wealth from age

wealth = poor

(prior = 0.760718)

| 1 | mean | cov |
|---|---|---|
| age | 37.374 | 198.935 |

density

0.025

0.015

0.005

20  30  40  50  60  70

age

wealth = rich

(prior = 0.239282)

| 1 | mean | cov |
|---|---|---|
| age | 44.7727 | 111.618 |

density

0.037

0.029

0.021

0.013

0.005

wealth values:  poor  rich

prob

1

0.6

0.2

*rich*

*Poor*

20  30  40  50  60  70  80  90

age

# Learning modelyear , mpg ---> maker

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

$$P(x|y=i) = \mathcal{N}(\mu_i, \Sigma_i)$$



maker = america
(prior = 0.625)

| 1 | mean | cov | |
|---|---|---|---|
| mpg | 20.0335 | 41.4785 | 15.2912 |
| modelyear | 75.5918 | 15.2912 | 13.3983 |

maker = asia
(prior = 0.201531)

| 1 | mean | cov | |
|---|---|---|---|
| mpg | 30.4506 | 37.0887 | 12.6427 |
| modelyear | 77.443 | 12.6427 | 13.3268 |

maker = europe
(prior = 0.173469)

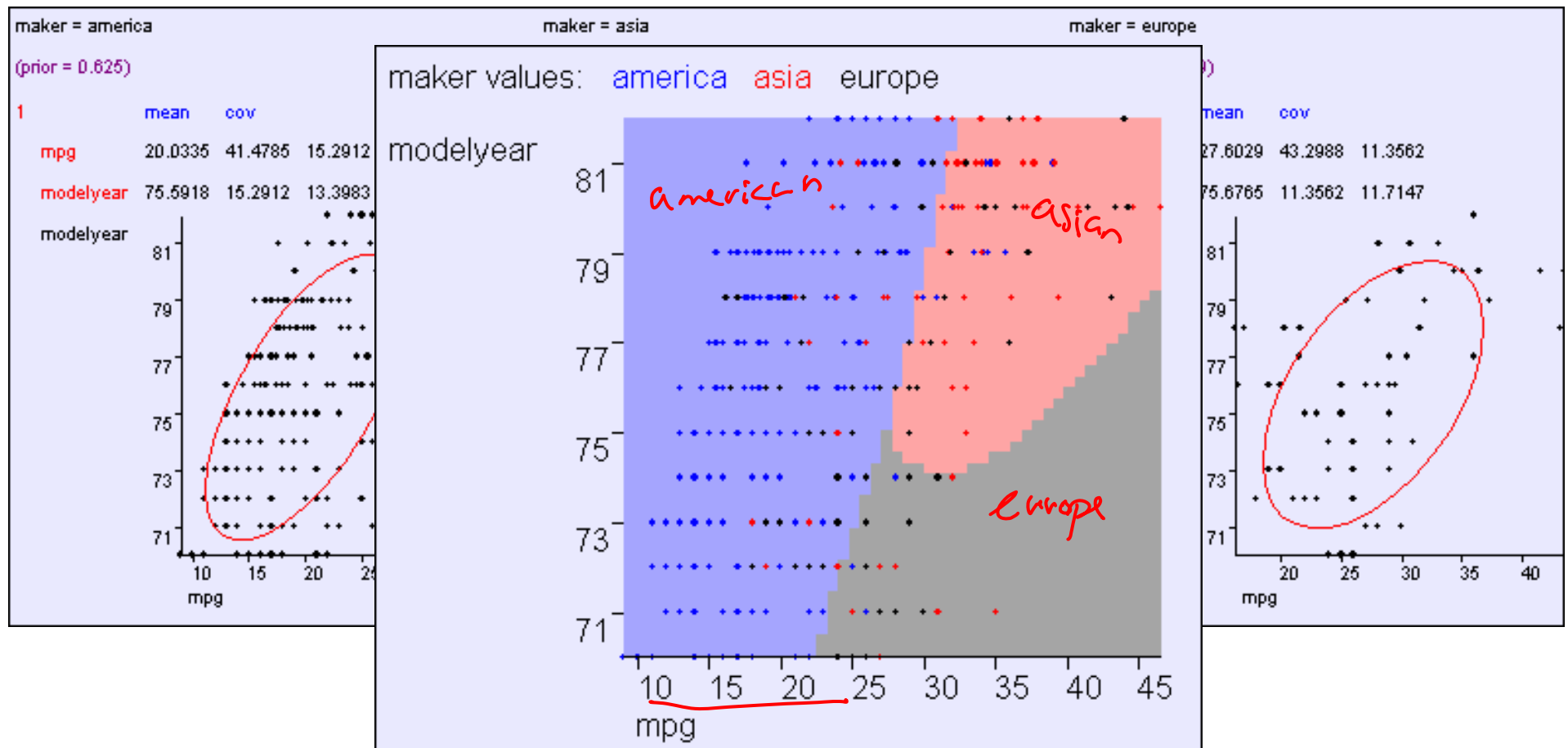| 1 | mean | cov | |
|---|---|---|---|
| mpg | 27.6029 | 43.2988 | 11.3562 |
| modelyear | 75.6765 | 11.3562 | 11.7147 |

# General: $O(m^2)$ parameters

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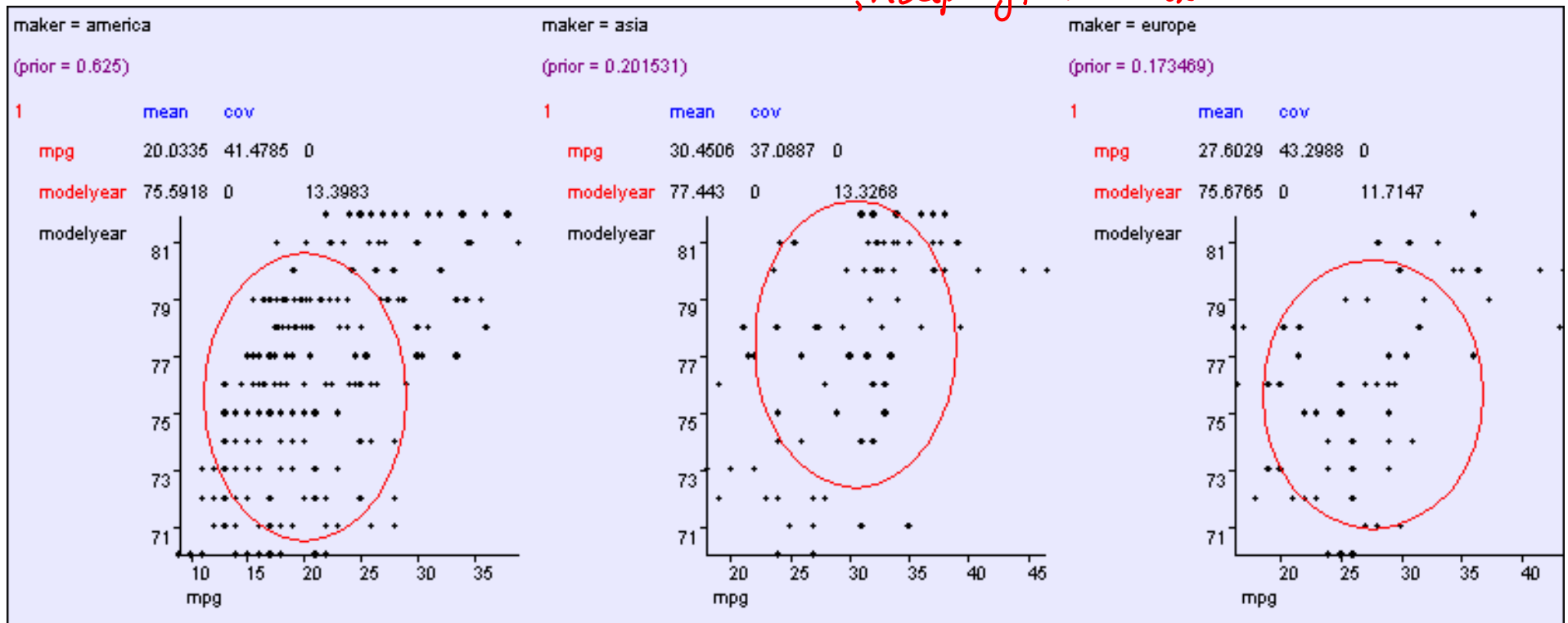
maker = america     maker = asia     maker = europe

(prior = 0.625)

1    mean    cov

mpg    20.0335   41.4785   15.2912

modelyear   75.5918   15.2912   13.3983

modelyear

maker values:   america   asia   europe

modelyear

american

asian

europe

mpg

27.6029   43.2988   11.3562

75.6765   11.3562   11.7147

# Aligned: *O(m)* parameters

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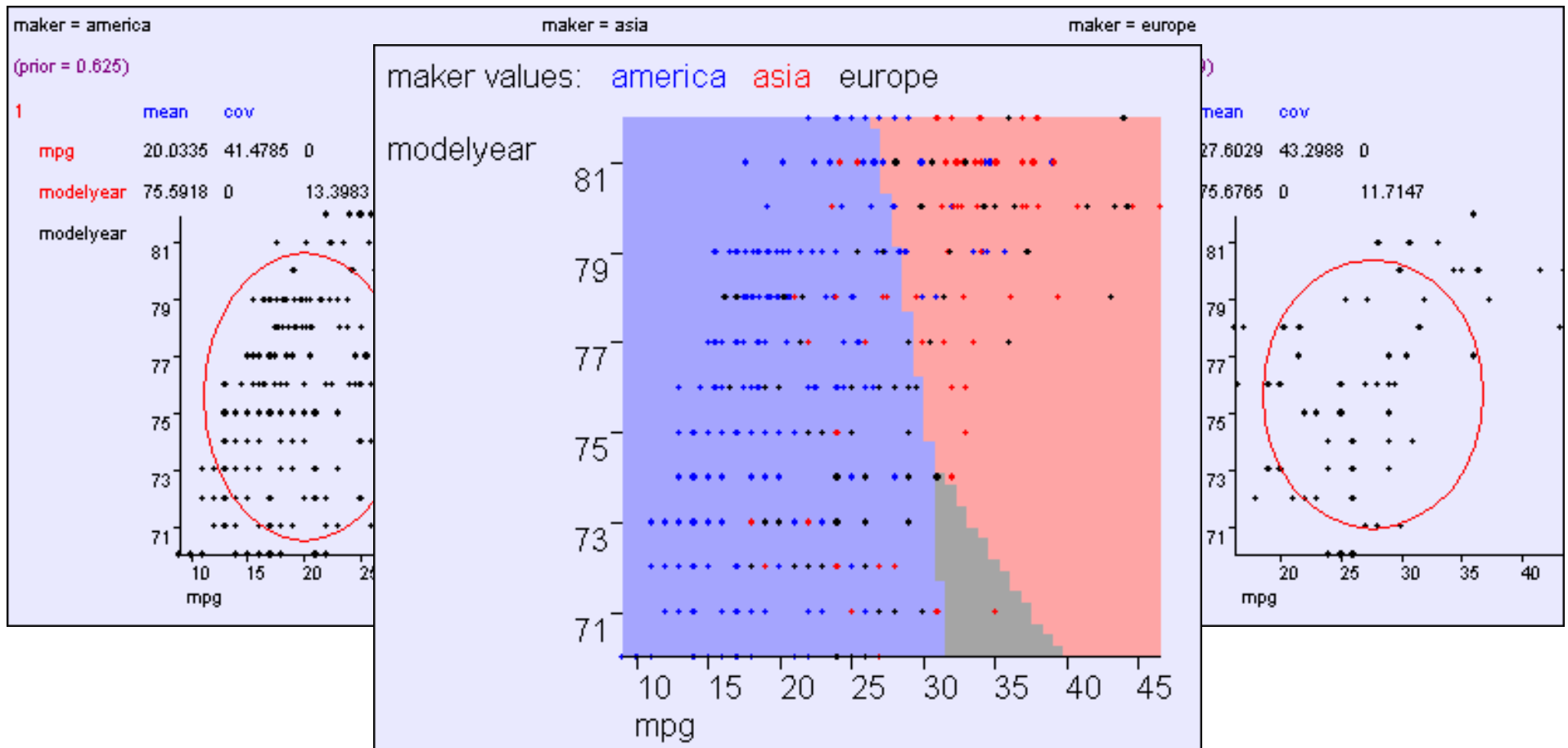
*features indep. given class*



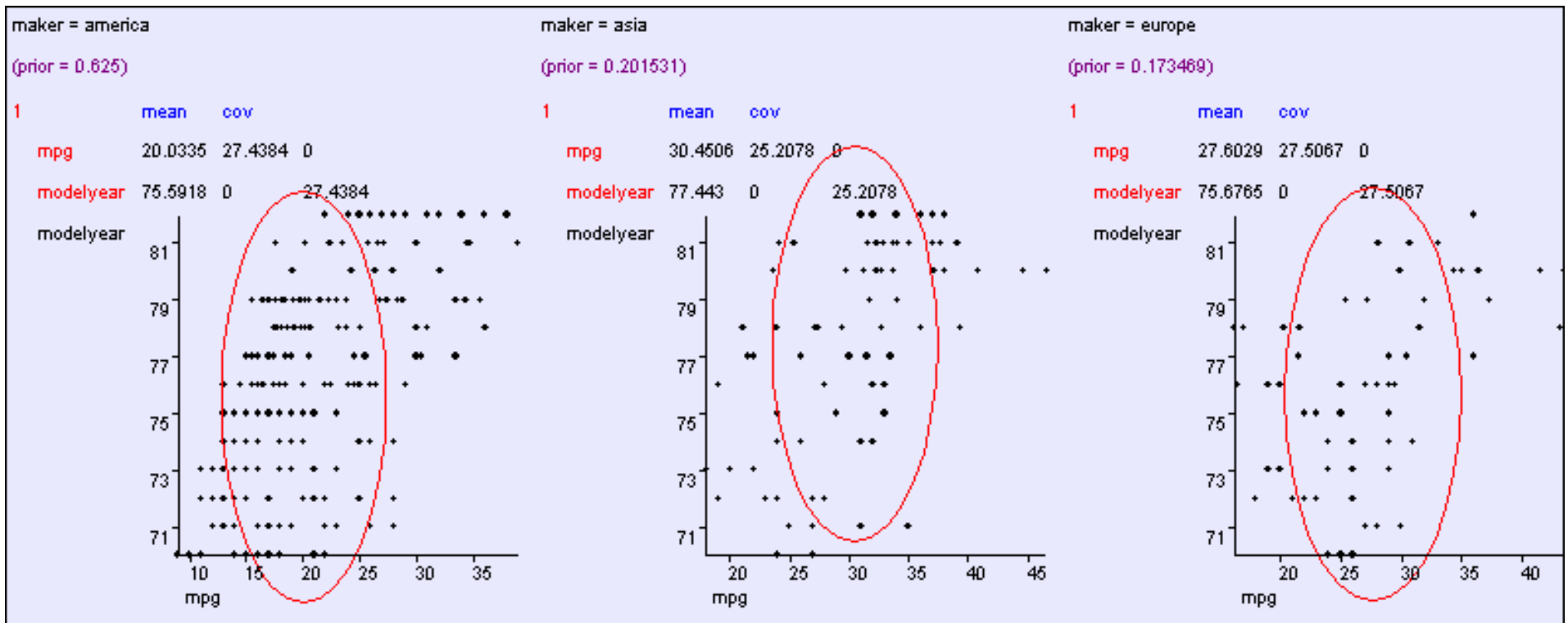| maker = america | | | maker = asia | | | maker = europe | | |
|---|---|---|---|---|---|---|---|---|
| (prior = 0.625) | | | (prior = 0.201531) | | | (prior = 0.173469) | | |
| 1 | mean | cov | 1 | mean | cov | 1 | mean | cov |
| mpg | 20.0335 | 41.4785  0 | mpg | 30.4506 | 37.0887  0 | mpg | 27.6029 | 43.2988  0 |
| modelyear | 75.5918 | 0   13.3983 | modelyear | 77.443 | 0   13.3268 | modelyear | 75.6765 | 0   11.7147 |

# Aligned: *O(m)* parameters

$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma^2_m \end{pmatrix}$$
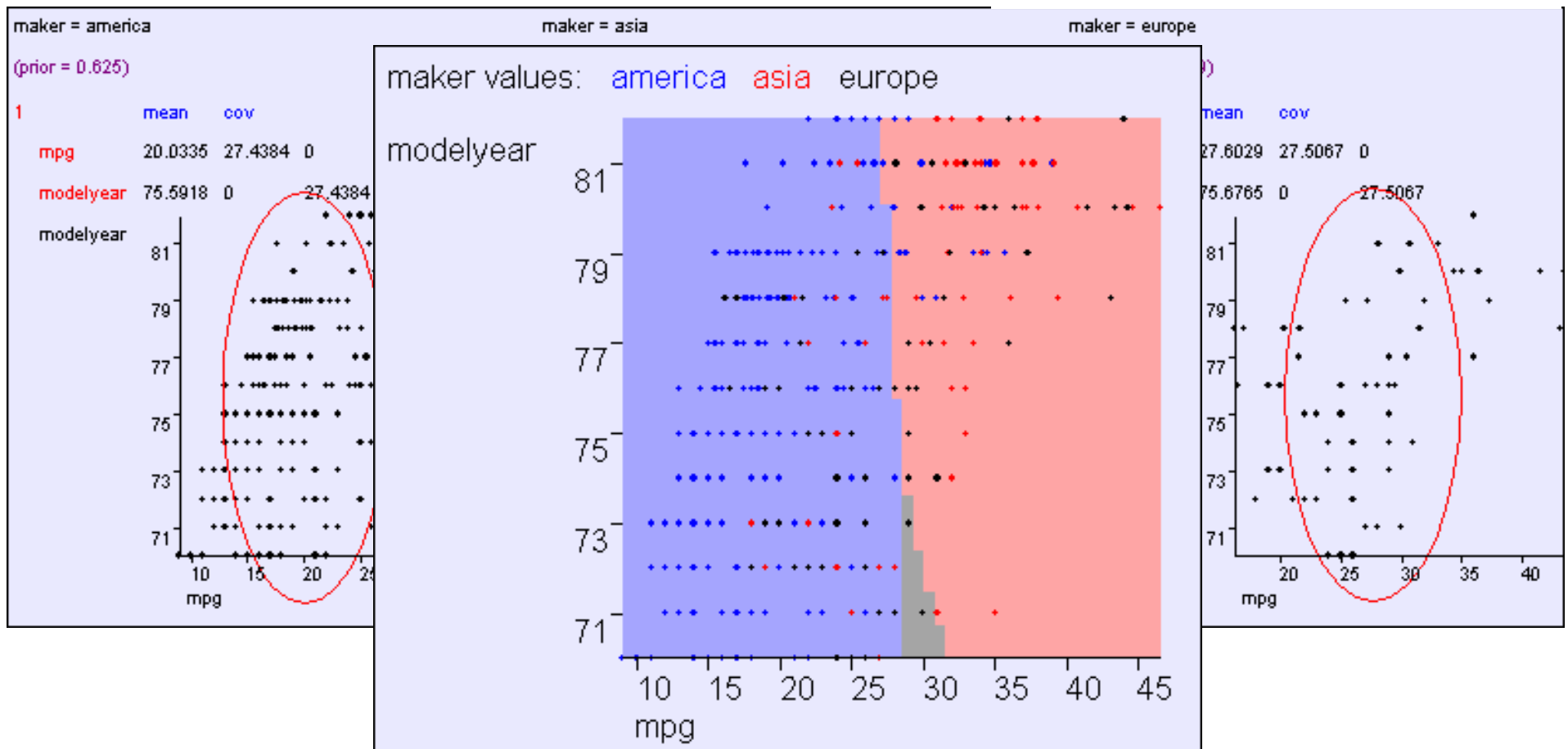
# Spherical: *O(1)* cov parameters

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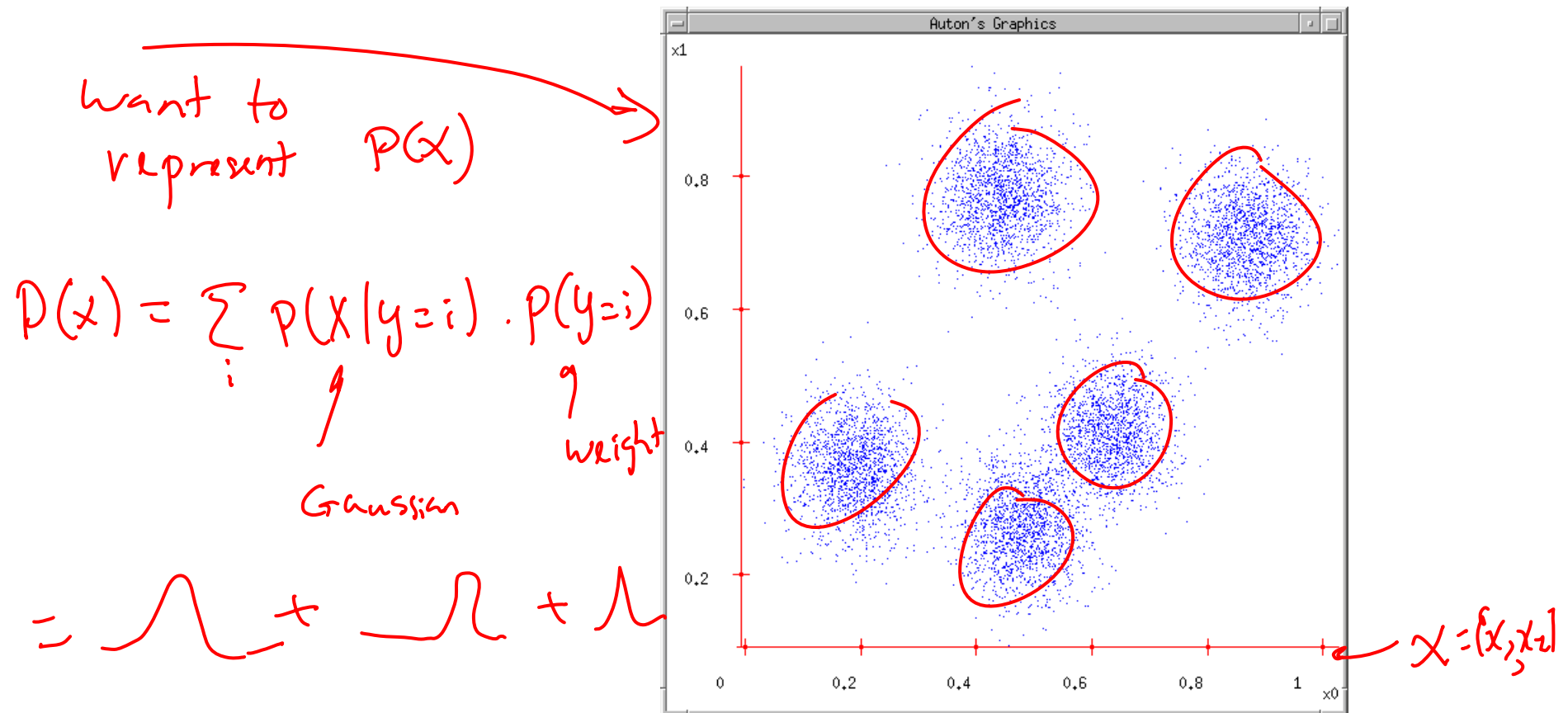
*every class same sigma*



maker = america (prior = 0.625)

| 1 | mean | cov | |
|---|---|---|---|
| mpg | 20.0335 | 27.4384 | 0 |
| modelyear | 75.5918 | 0 | 27.4384 |

maker = asia (prior = 0.201531)

| 1 | mean | cov | |
|---|---|---|---|
| mpg | 30.4506 | 25.2078 | 0 |
| modelyear | 77.443 | 0 | 25.2078 |

maker = europe (prior = 0.173469)

| 1 | mean | cov | |
|---|---|---|---|
| mpg | 27.6029 | 27.5067 | 0 |
| modelyear | 75.6765 | 0 | 27.5067 |

# Spherical: *O(1)* cov parameters

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

# Next… back to Density Estimation

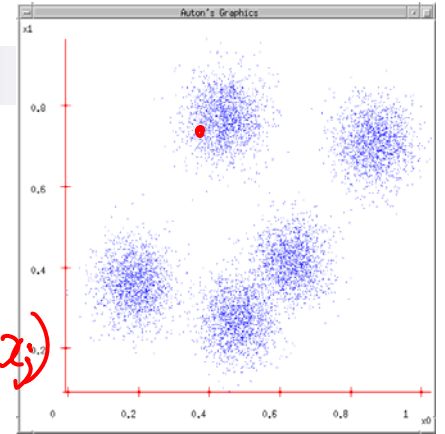## What if we want to do density estimation with multimodal or clumpy data?

Want to represent $P(X)$

$$P(x) = \sum_i P(X | y=i) \cdot P(y=i)$$

Gaussian

weight

$$= \underset{}{\bigwedge} + \underset{}{\bigwedge} + \bigwedge$$

$X = [x_1, x_2]$

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



- **MLE:**

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

  $= \text{argmax } \log \prod_j P(y_j, x_j) = \text{argmax} \sum_j \log P(y_j, x_j)$

  *Same as MLE*

  *typically easy*

- **But we don't know $y_j$'s!!!**

- **Maximize marginal likelihood:**

  - argmax $\prod_j P(x_j) = \text{argmax} \prod_j \sum_{i=1}^k P(y_j=i, x_j)$

  $= \text{argmax} \sum_j \log \sum_{i=1}^k P(y_j=i, x_j)$

  *log sum → typically hard*

# Special case: spherical Gaussians and hard assignments

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

$P(y=i, x_j)$    8 features    $\Sigma = \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}$

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

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

"=" $\frac{1}{k}$ ≈ uniform    all clusters about same size

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

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

$$= \sum_{j=1}^{m} -\frac{1}{2\sigma^2} \| x_j - \mu_{C(j)} \|^2$$

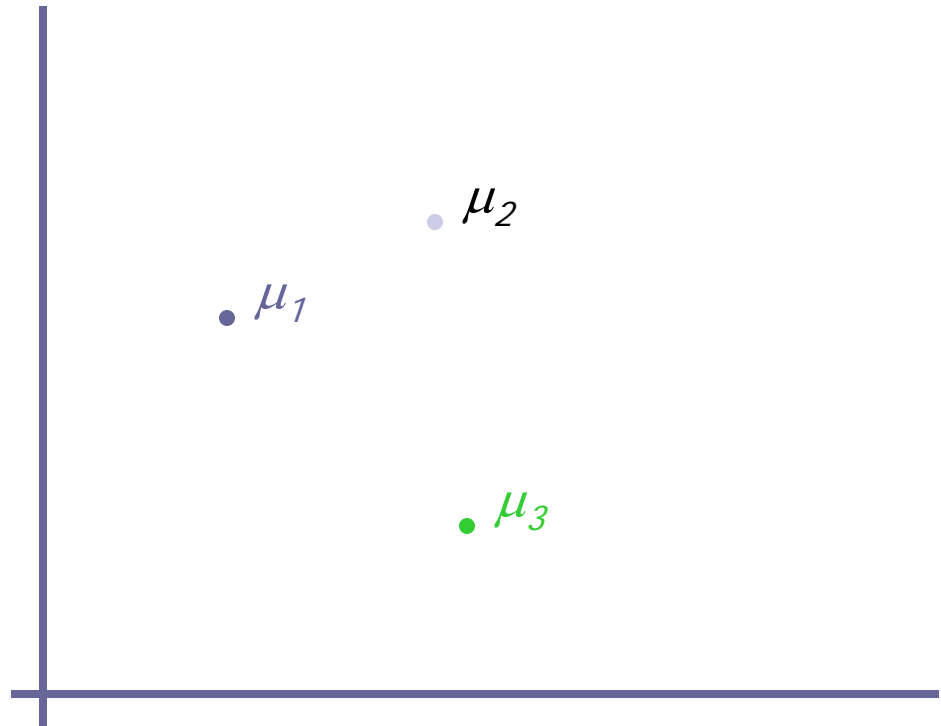- Same as K-means!!!

narrow special case of fitting mixtures of Gaussians

# The GMM assumption

- There are <u>k components</u>
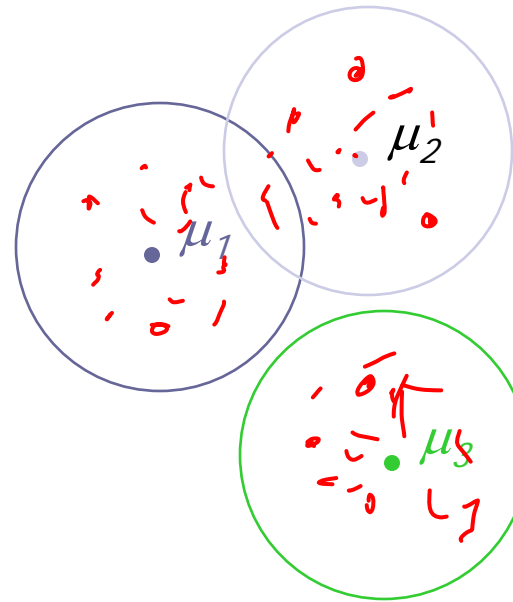
- Component $i$ has an associated mean vector $\mu_i$

$\bullet\ \mu_2$

$\bullet\ \mu_1$

$\bullet\ \mu_3$

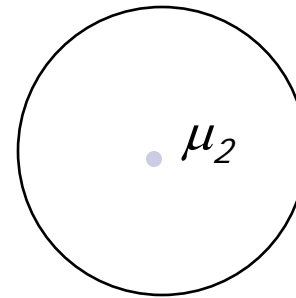# The GMM assumption

- 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 $\sigma^2 I$   e.g., $\Sigma = \begin{pmatrix} \sigma^2 & 0 \\ 0 & \ddots \\ & & \sigma^2 \end{pmatrix}$

  Each data point is generated according to the following recipe:

# 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:

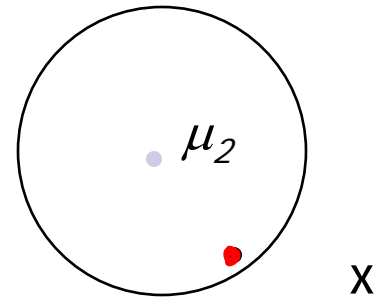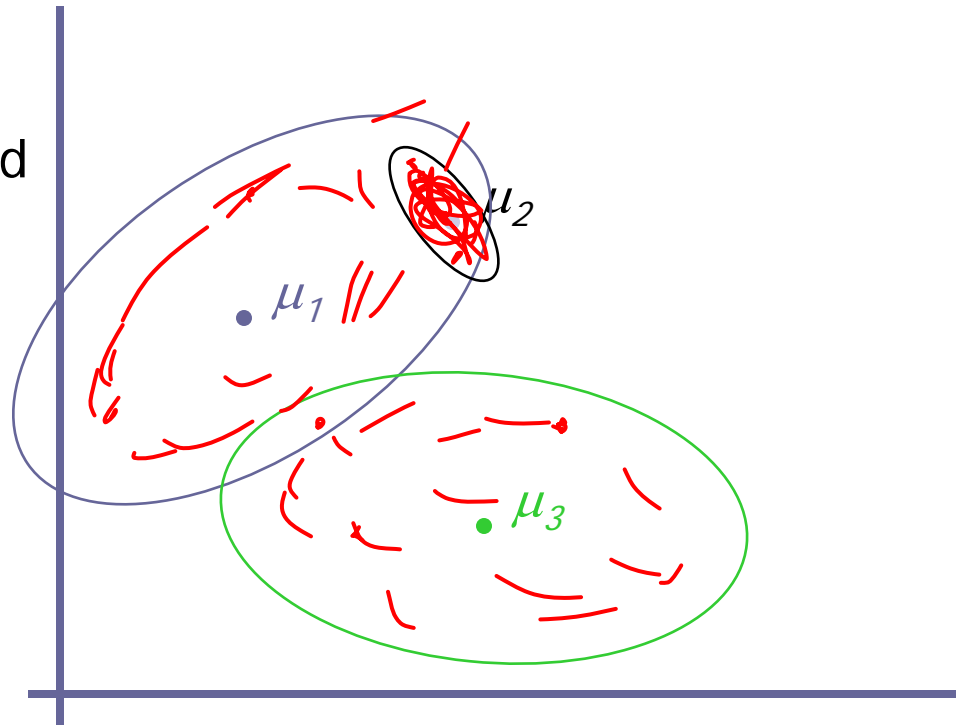1. Pick a component at random: Choose component i with probability $P(y=i)$

$\mu_2$

# 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 $\sigma^2 I$

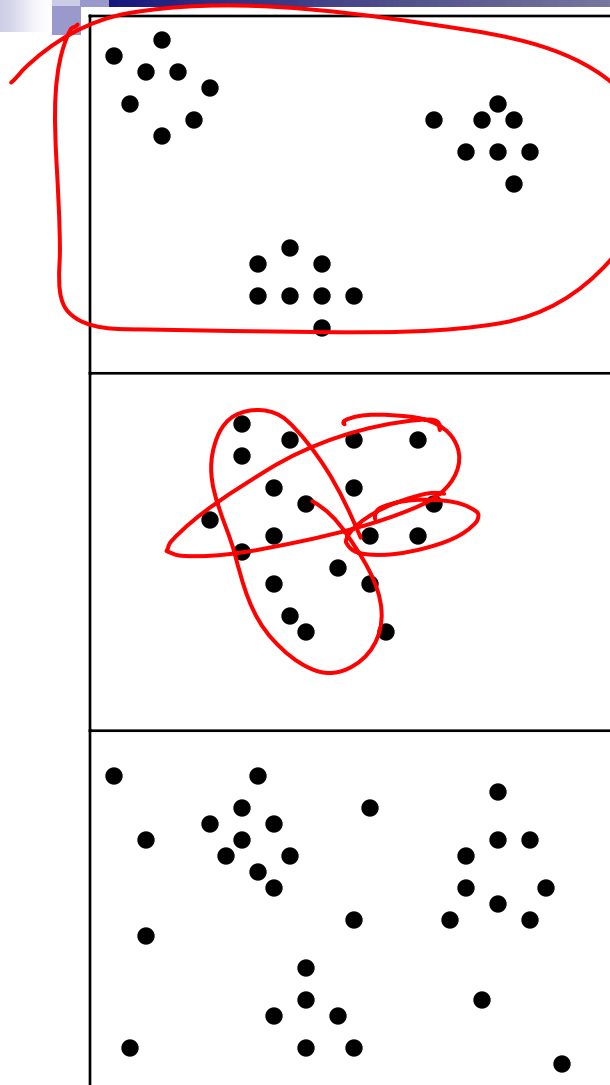Each data point is generated according to the following recipe:

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

2. Datapoint ~ $N(\mu_{j'}, \sigma^2 I)$



$\mu_2$

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 ~ $N(\mu_i, \Sigma_i)$

general covariance matrix

# Unsupervised Learning: not as hard as it looks

*well-separated*

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

# 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:

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

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

*Handwritten annotations (red):*

log

$= \sum_{j=1}^{m} \log \sum_{i=1}^{k}$

defn. of GMM

$x_j \leftarrow$ Observed

assumption: $x_j \sim$ GMM

$P(x_j) = \sum_i P(y=i) \cdot P(x_j \mid y_{=i})$

don't observe $y_j$ $\Rightarrow$ max $P(x_j)$

Gaussian

# Special case 2: spherical Gaussians and soft assignments

$$\to \Sigma_i = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

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

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

*addresses overlap*

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

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

*P(x,y=i)*

*guess* $\mu_1 \dots \mu_k$ $\to$ *compute* $P(x_j | y_i)$

*recompute centers weighed by prob.*

# 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.

*for each $x$*

$$P(x|y) = N(\mu, \sigma^2)$$
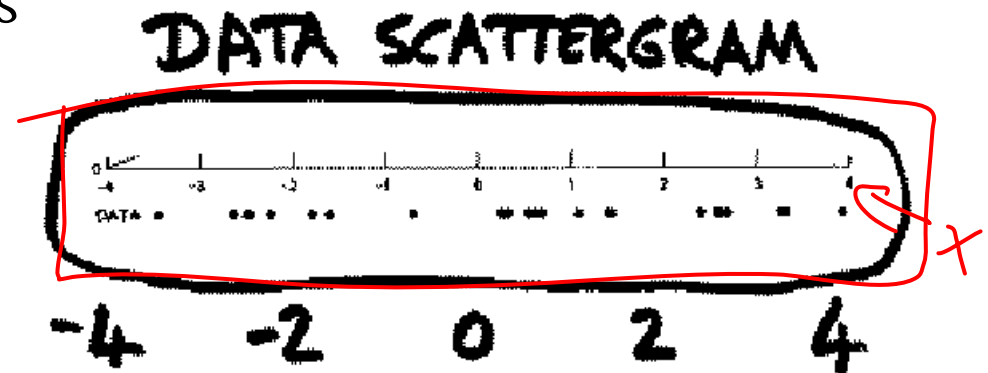
*given $\mu$*

Suppose $x$'s are 1-dimensional.

**(From Duda and Hart)**

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

*Set $\sigma = 1$*

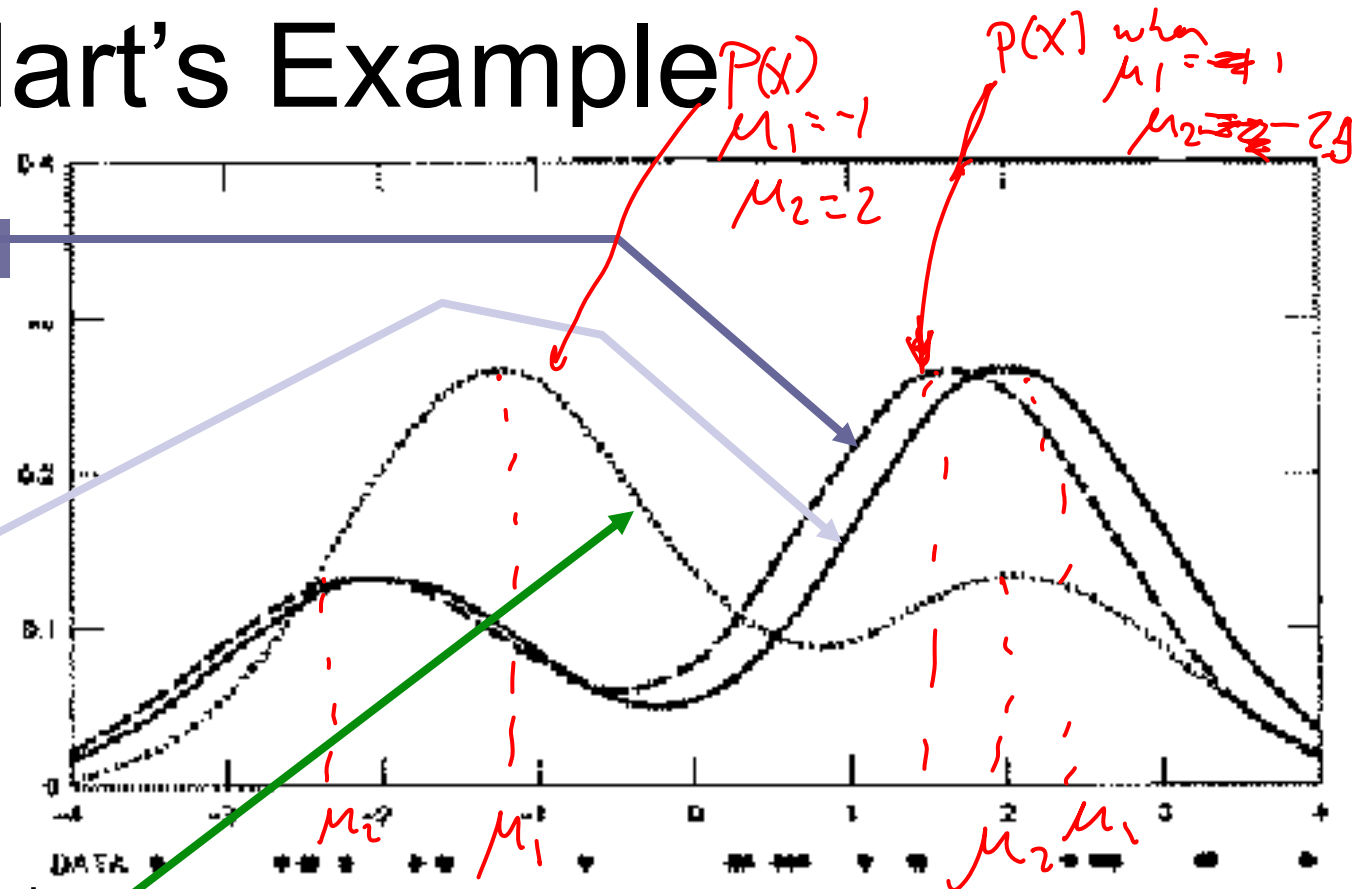$P(y_1) = 1/3$     $P(y_2) = 2/3$     $\sigma = 1$

There are 25 unlabeled datapoints

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

**DATA SCATTERGRAM**

# Duda & Hart's Example



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

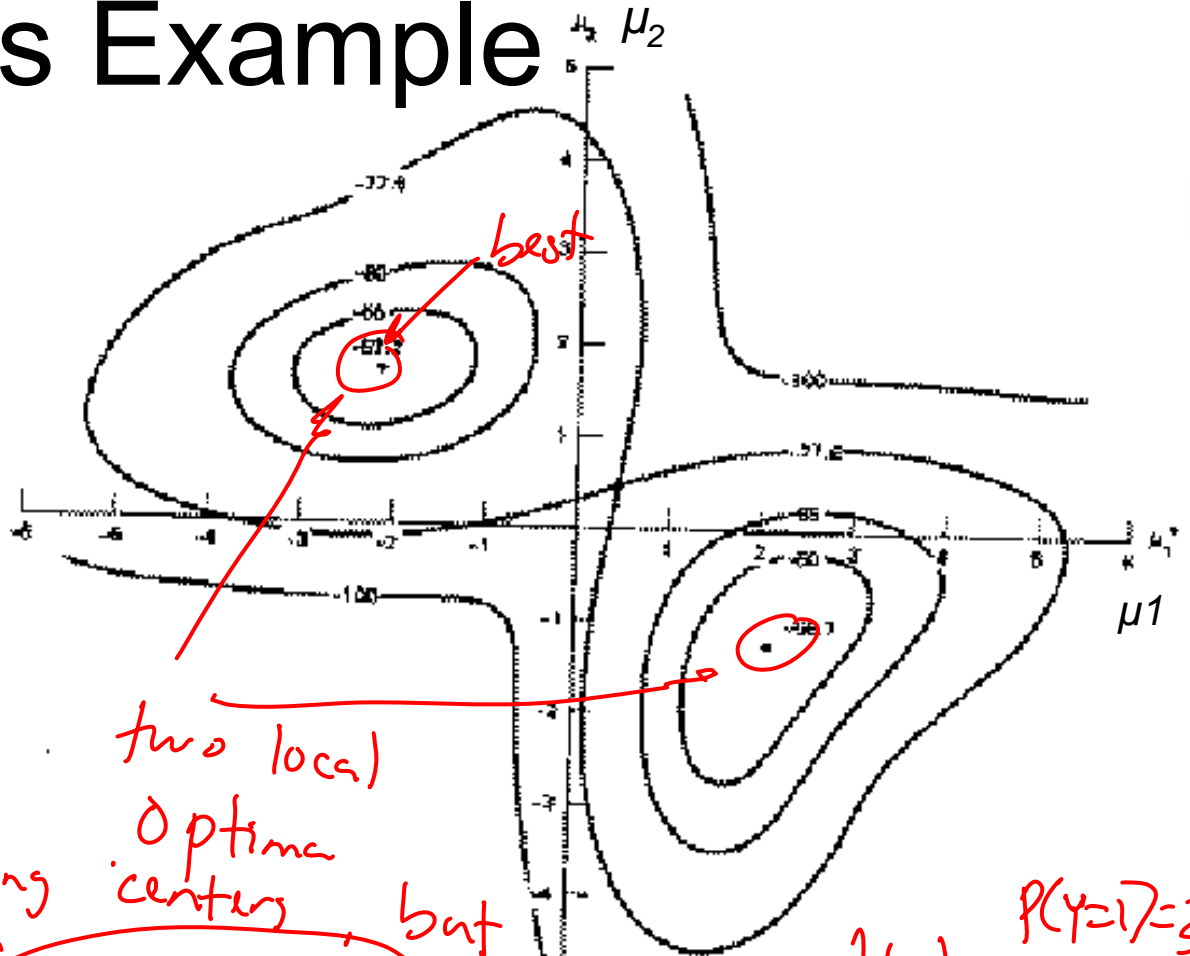We can also graph the true function from which the data was randomly generated.

Handwritten annotations on figure:
$P(x)$ $\mu_1 = -1$ $\mu_2 = 2$
$P(x)$ when $\mu_1 = +1$, $\mu_2 = -2$
$\mu_2$? $\mu_1$ $\mu_2$ $\mu_1$

- They are close.  Good.

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

  $P(y=1) = \frac{2}{3}$    $P(y=2) = \frac{1}{3}$

- 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

max $P(X|\mu)$
$\mu$

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

*best*

*two local*
*Optima*
*flipping centers, but we knew that*

$P(y=1)=\frac{2}{3}$
$P(y=2)=\frac{1}{3}$

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

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

&ast; corresponds to switching $y_1$ with $y_2$.

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

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

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

- The normal max likelihood trick:

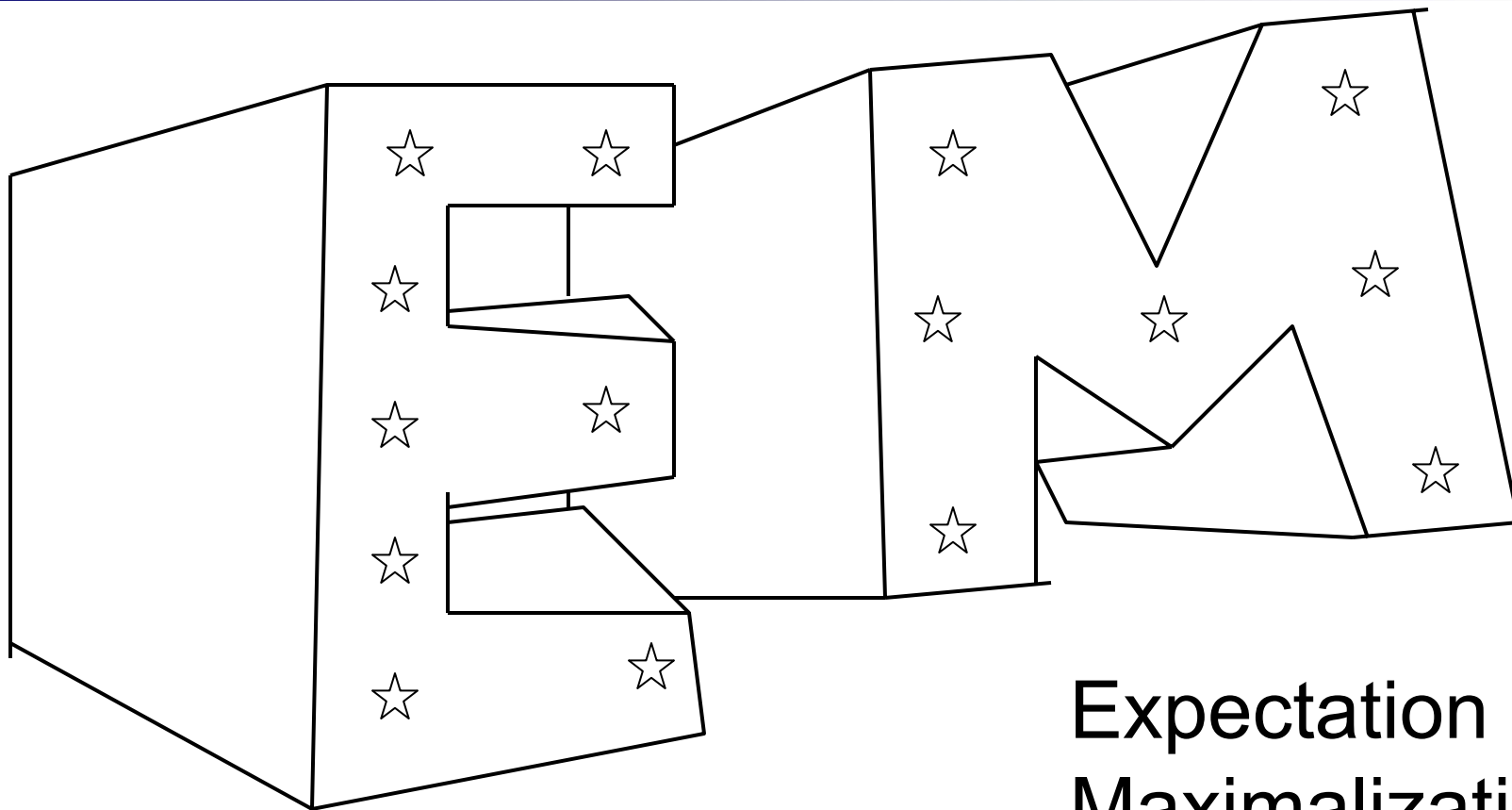  Set $\frac{\partial}{\partial \mu_i}$ log Prob (....) = 0

  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…

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 inference of Hidden Markov Models (future lecture).

# 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

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

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

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

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

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

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 $\mu$ now?

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

# 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 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

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

©2005-2007 Carlos Guestrin

# 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

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

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

**E-step**

$$\mu^{(t+1)} = \frac{b^{(t)}+c}{6\left(b^{(t)}+c+d\right)}$$

**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.**

# 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

Remember:

We have unlabeled data $x_1$ $x_2$ … $x_m$

We know there are k classes

We know $P(y_1)$ $P(y_2)$ $P(y_3)$ … $P(y_k)$

We <u>don't</u> know $\mu_1$ $\mu_2$ .. $\mu_k$

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

$$= \mathrm{p}\left(x_1...x_m \middle| \mu_1...\mu_k\right)$$

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

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

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

# 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 $P(y=i|x_j)$:

$$\mu_i = \frac{\sum\limits_{j=1}^{m} P(y=i|x_j) x_j}{\sum\limits_{j=1}^{m} P(y=i|x_j)}$$

# E.M. for GMMs

**E-step**

Compute "expected" classes of all datapoints for each class

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

*Just evaluate a Gaussian at $x_j$*

**M-step**

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

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

©2005-2007 Carlos Guestrin

# E.M. Convergence

- EM is coordinate ascent on an interesting potential function

- Coord. ascent for bounded pot. func. $\rightarrow$ convergence to a local optimum guaranteed

- 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

# E.M. for General GMMs

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

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

**E-step**

Compute "expected" classes of all datapoints for each class

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

*Just evaluate a Gaussian at $x_j$*

M-step

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

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

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

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

$m$ = #records

# Gaussian Mixture Example: Start



p=0.333

p=0.333

0.333

# After first iteration

# After 2nd iteration



p=0.374

p=0.306

p=0.320

# After 3rd iteration



p=0.345

p=0.307

# After 4th iteration



p=0.331

p=0.288

# After 5th iteration



p=0.322

p=0.285

# 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



©2005-2007 Carlos Guestrin

# Resulting Density Estimator

# Three classes of assay

(each learned with it's own mixture model)



Compound =
IL-1
TNF
none

nucleus

# Resulting Bayes Classifier



Compound =
- IL-1
- TNF
- none
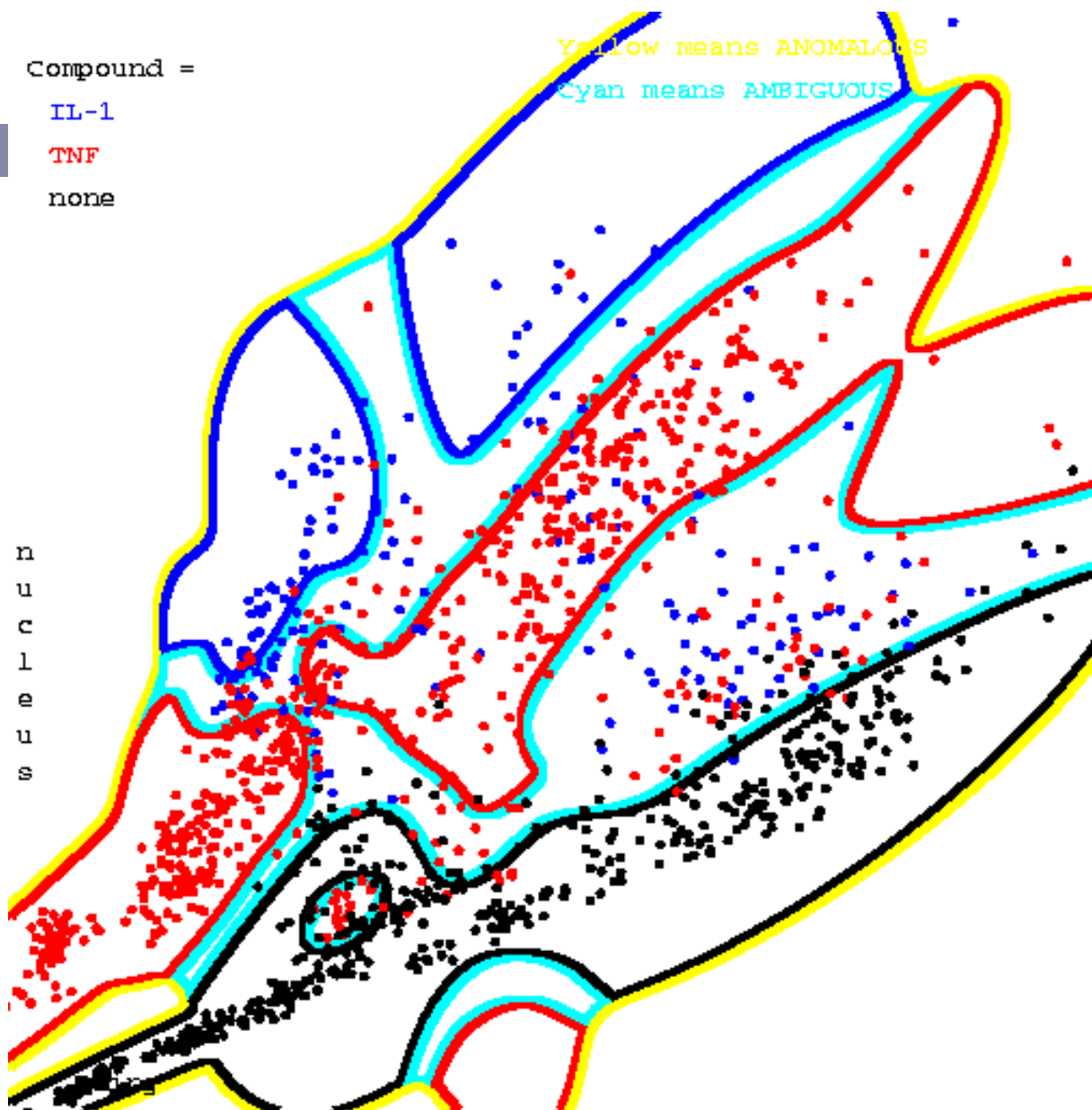
Resulting Bayes Classifier, using posterior probabilities to alert about ambiguity and anomalousness

**Yellow means anomalous**

**Cyan means ambiguous**

Compound =
IL-1
TNF
none

nucleus

Yellow means ANOMALOUS
Cyan means AMBIGUOUS

©2005-2007 Carlos Guestrin

# 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

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

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

# 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