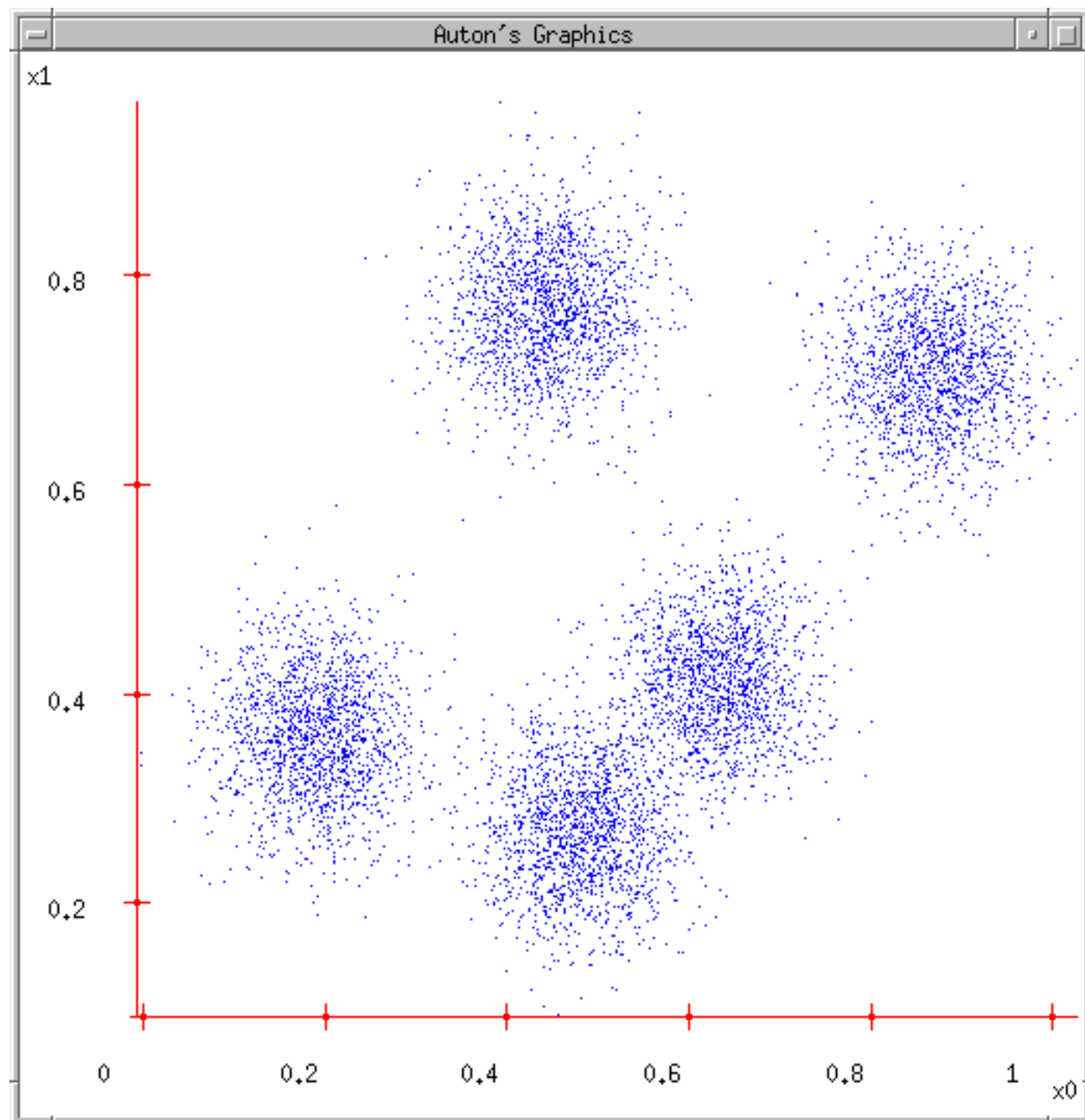


# CSE 575: Statistical Machine Learning

Jingrui He  
CIDSE, ASU

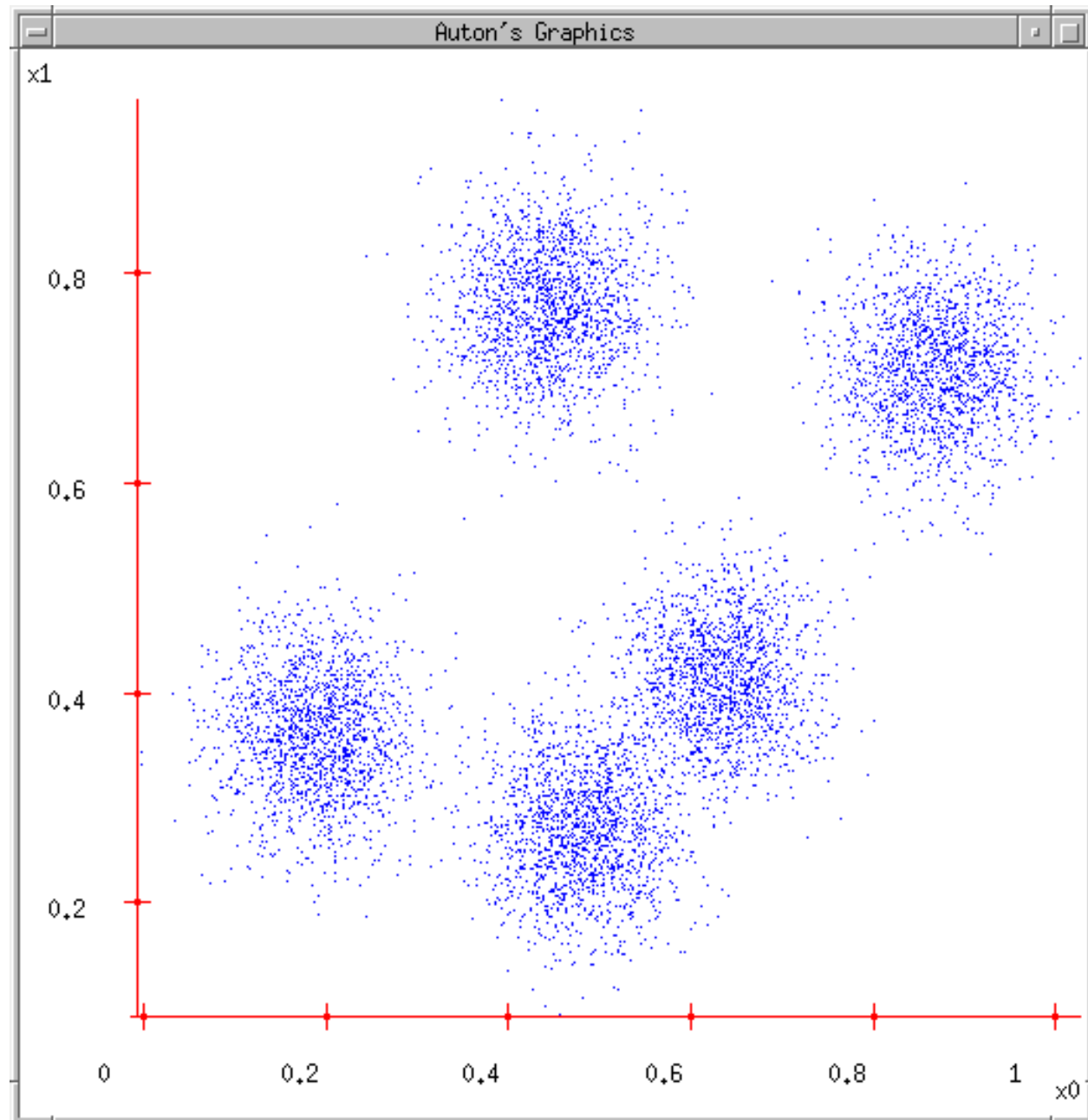
K-means, GMM, EM

# Some Data



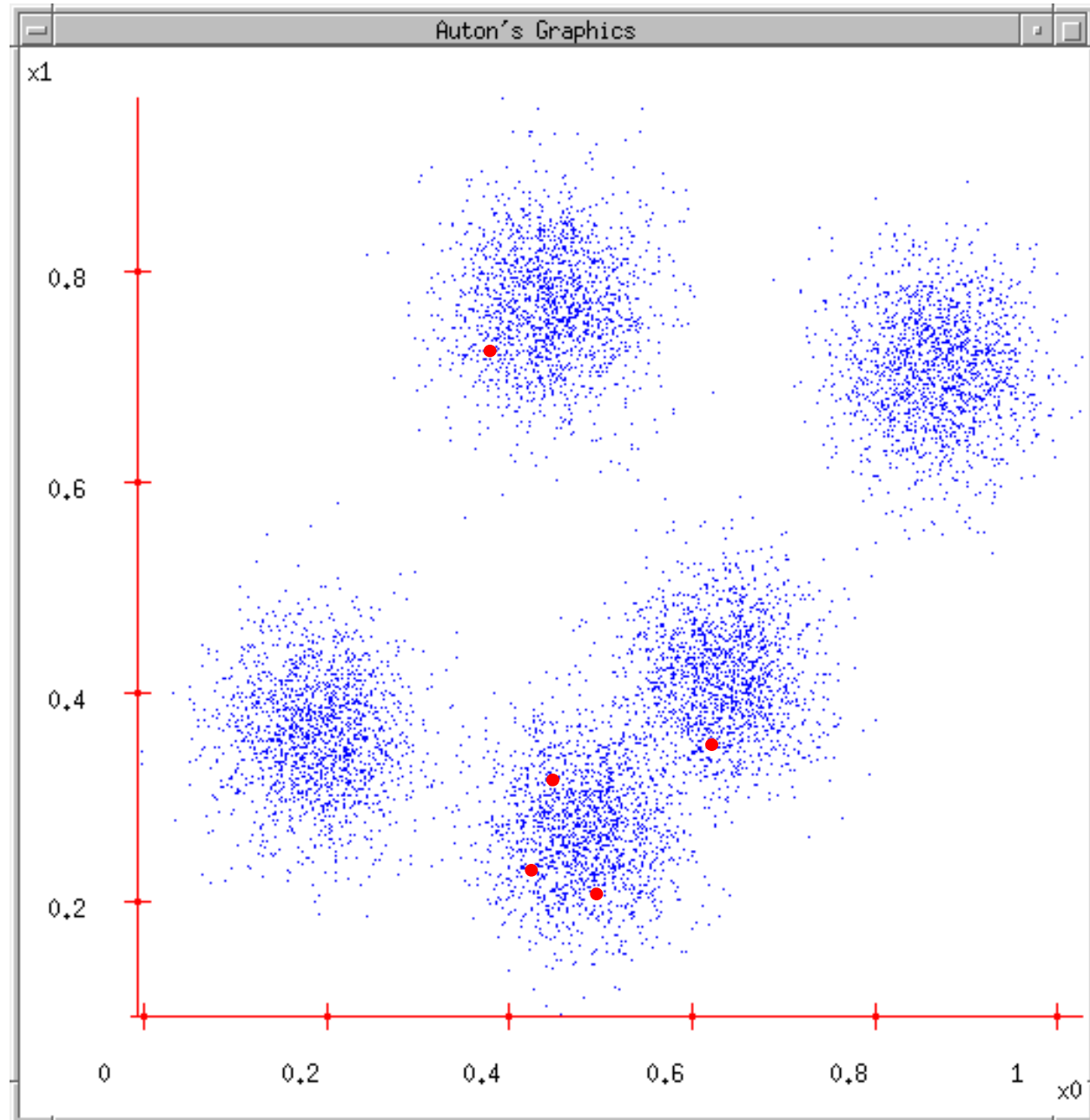
# K-means

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



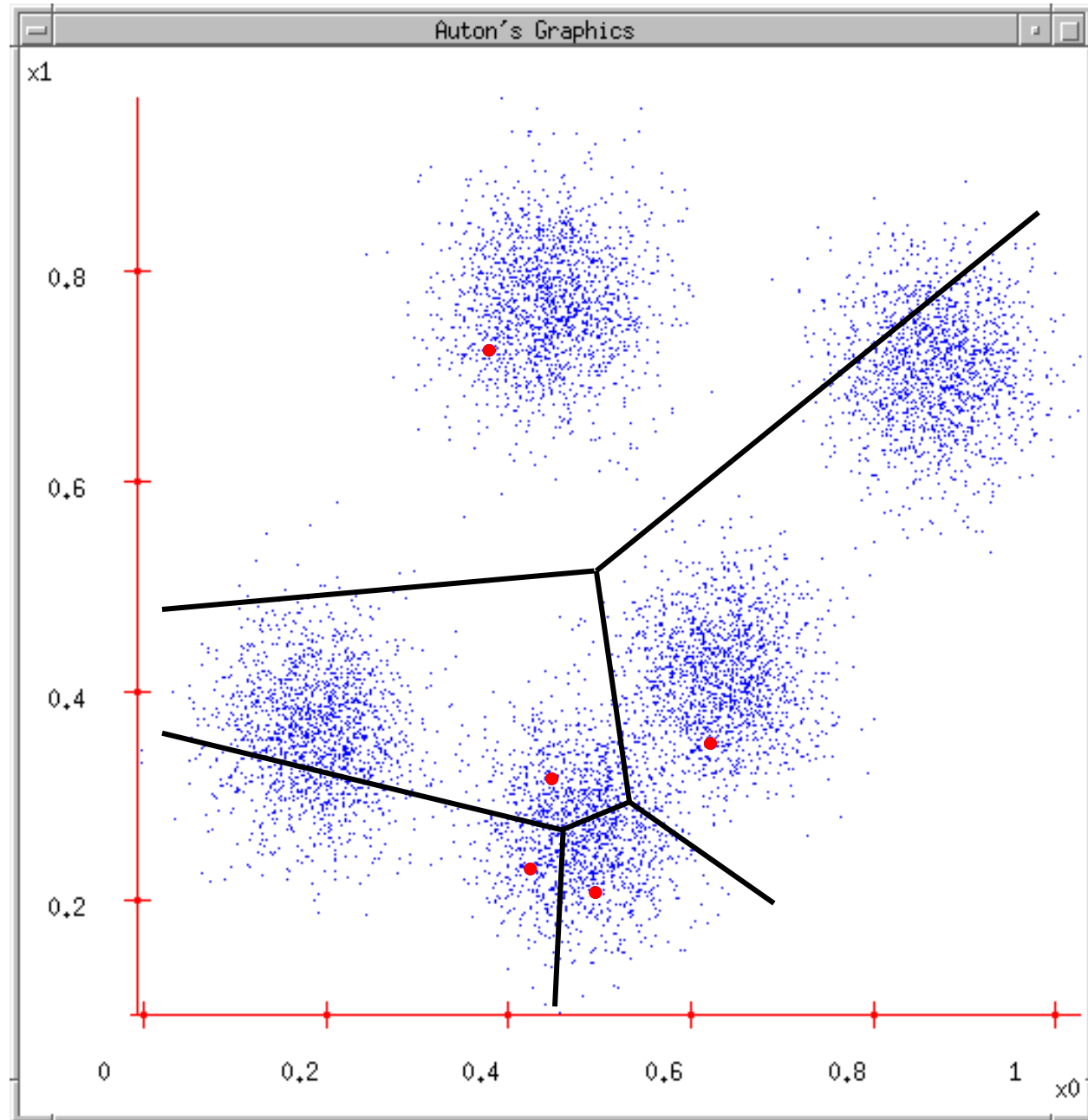
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



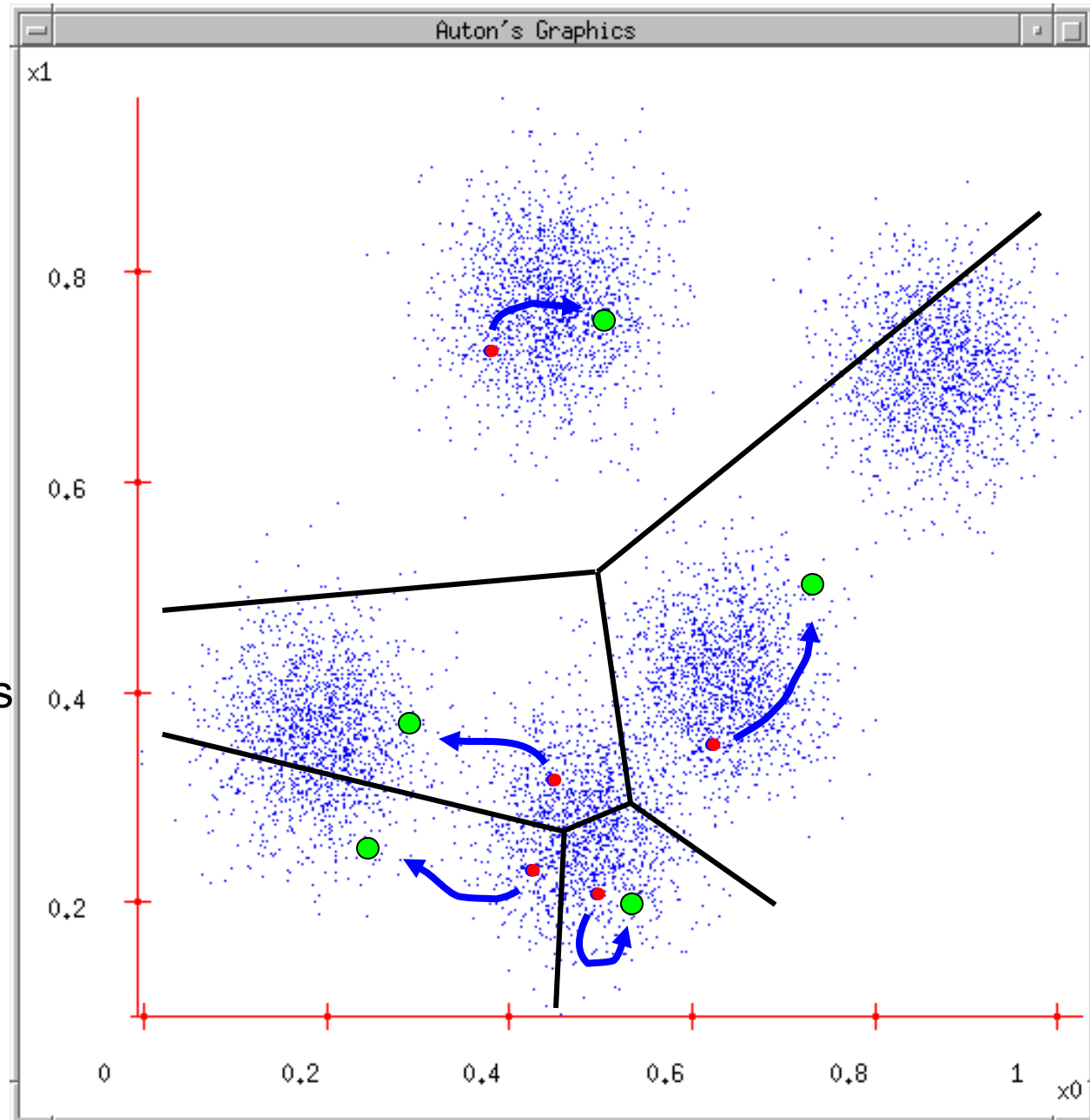
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each data point finds out which Center it's closest to  
(Thus each Center "owns" a set of data points)



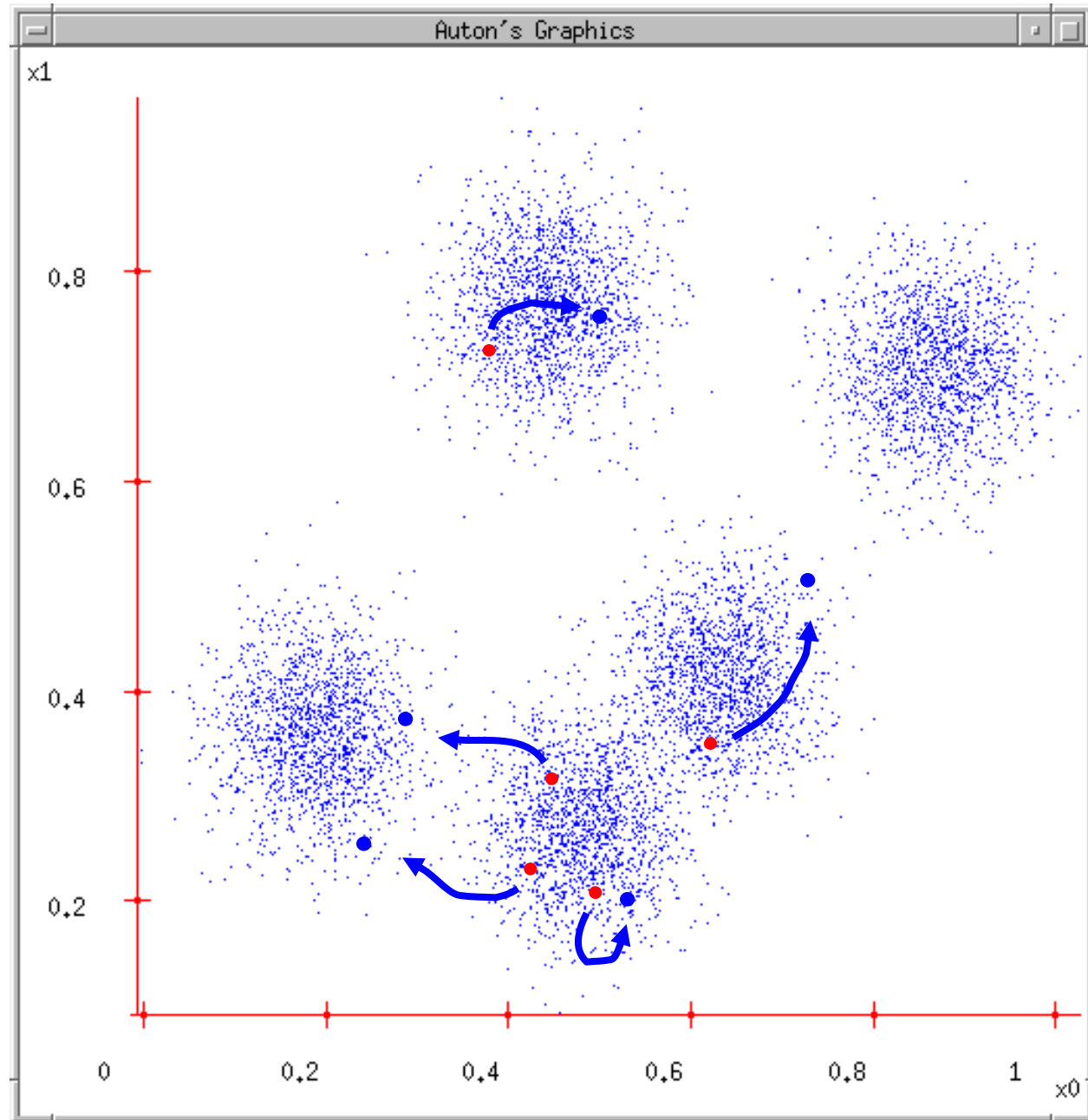
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each data point finds out which Center it's closest to
4. Each Center finds the centroid of the points it owns



# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each data point 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!





# K-means

- Randomly initialize  $k$  centers
  - $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$
- **Classify:** Assign each point  $j \in \{1, \dots, m\}$  to nearest center:
  - $C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$
- **Recenter:**  $\mu_i$  becomes centroid of its point:
  - $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C(j)=i} \|\mu - x_j\|^2$
  - Equivalent to  $\mu_i \leftarrow$  average of its points!

# What is K-means optimizing?

- Potential function  $F(\mu, C)$  of centers  $\mu$  and point allocations  $C$ :

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

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

# Does K-means converge??? Part 1

- Optimize potential function:

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

- Fix  $\mu$ , optimize C

# 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  $\mu$

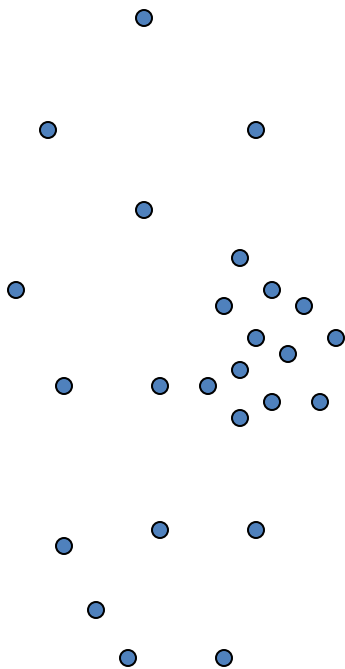
# Coordinate descent algorithms

- 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
- K-means is a coordinate descent algorithm!

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

# (One) bad case for k-means

- Clusters may overlap
- Some clusters may be “wider” than others



# Gaussian Bayes classifier reminder

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

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

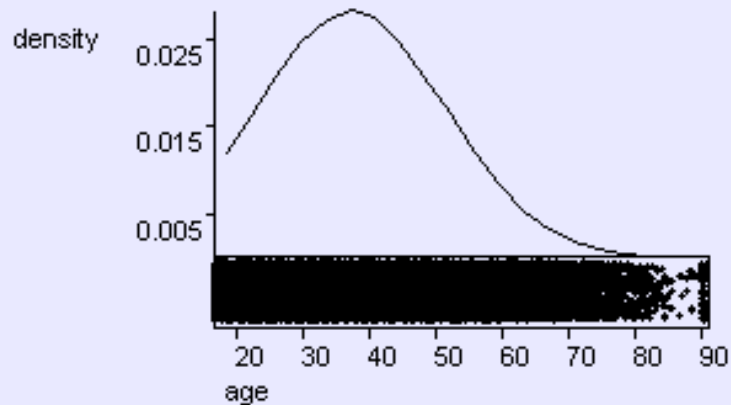
# Predicting wealth from age

wealth = poor

(prior = 0.760718)

1      mean    cov

age 37.374 198.935

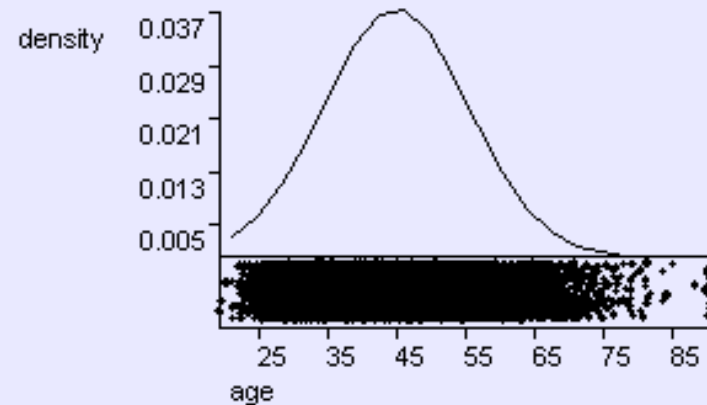


wealth = rich

(prior = 0.239282)

1      mean    cov

age 44.7727 111.618





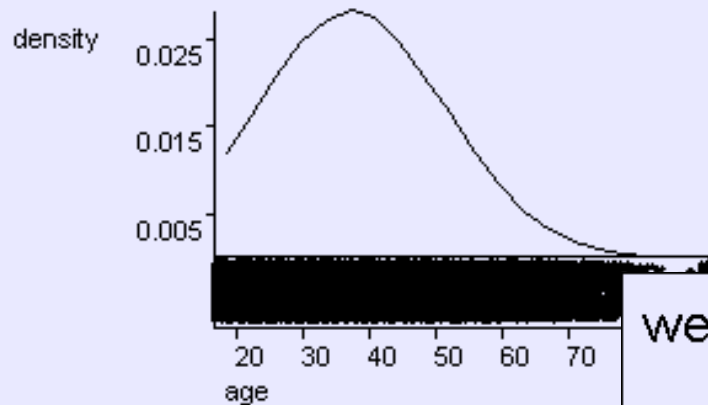
# Predicting wealth from age

wealth = poor

(prior = 0.760718)

1      mean    cov

age 37.374 198.935

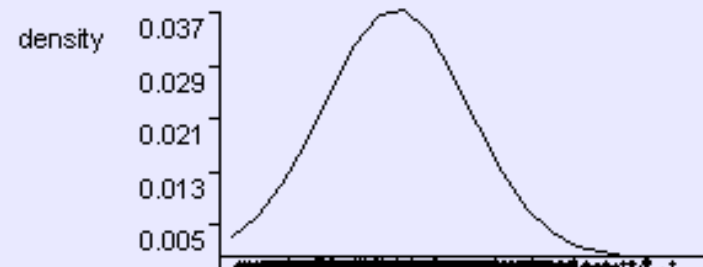


wealth = rich

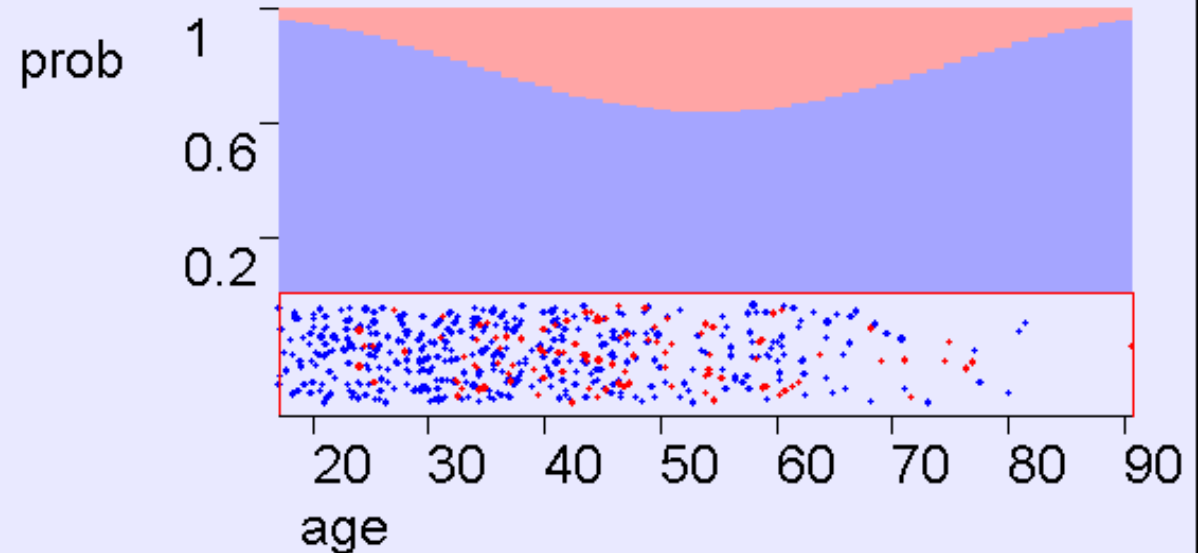
(prior = 0.239282)

1      mean    cov

age 44.7727 111.618

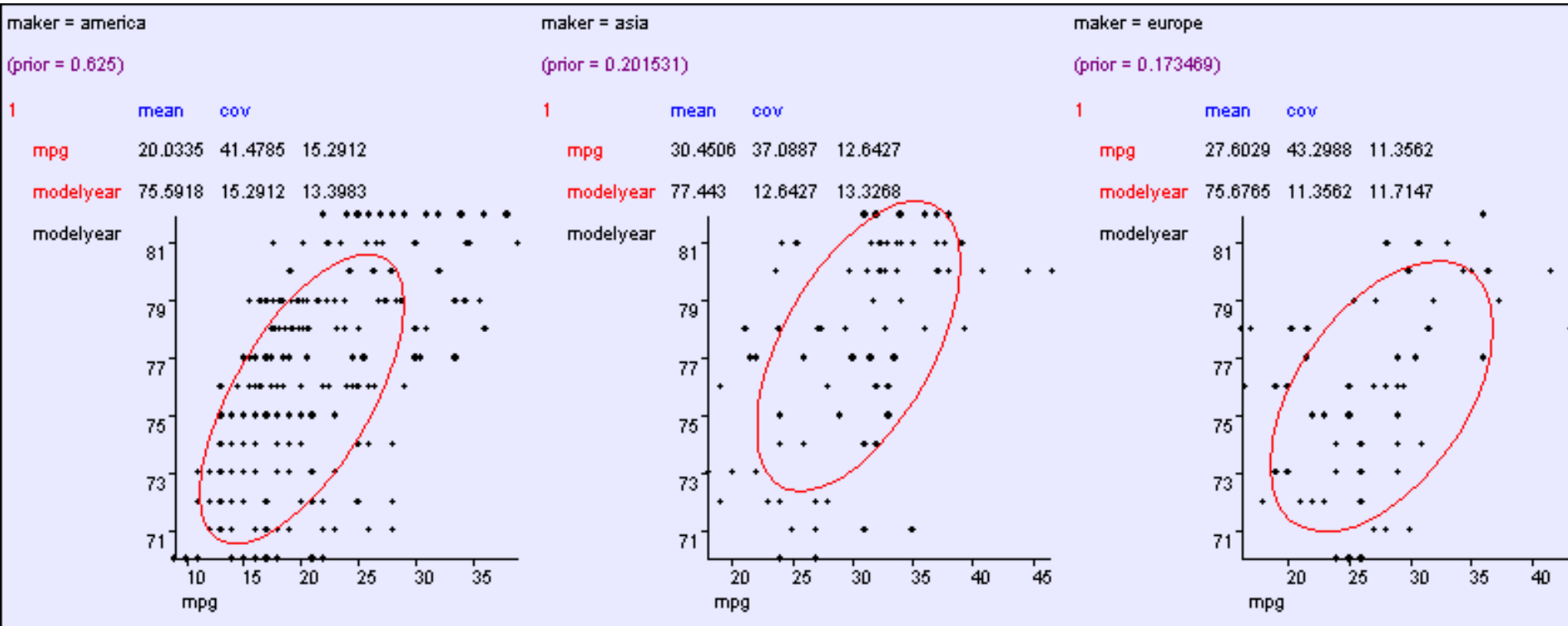


wealth values:    poor    rich



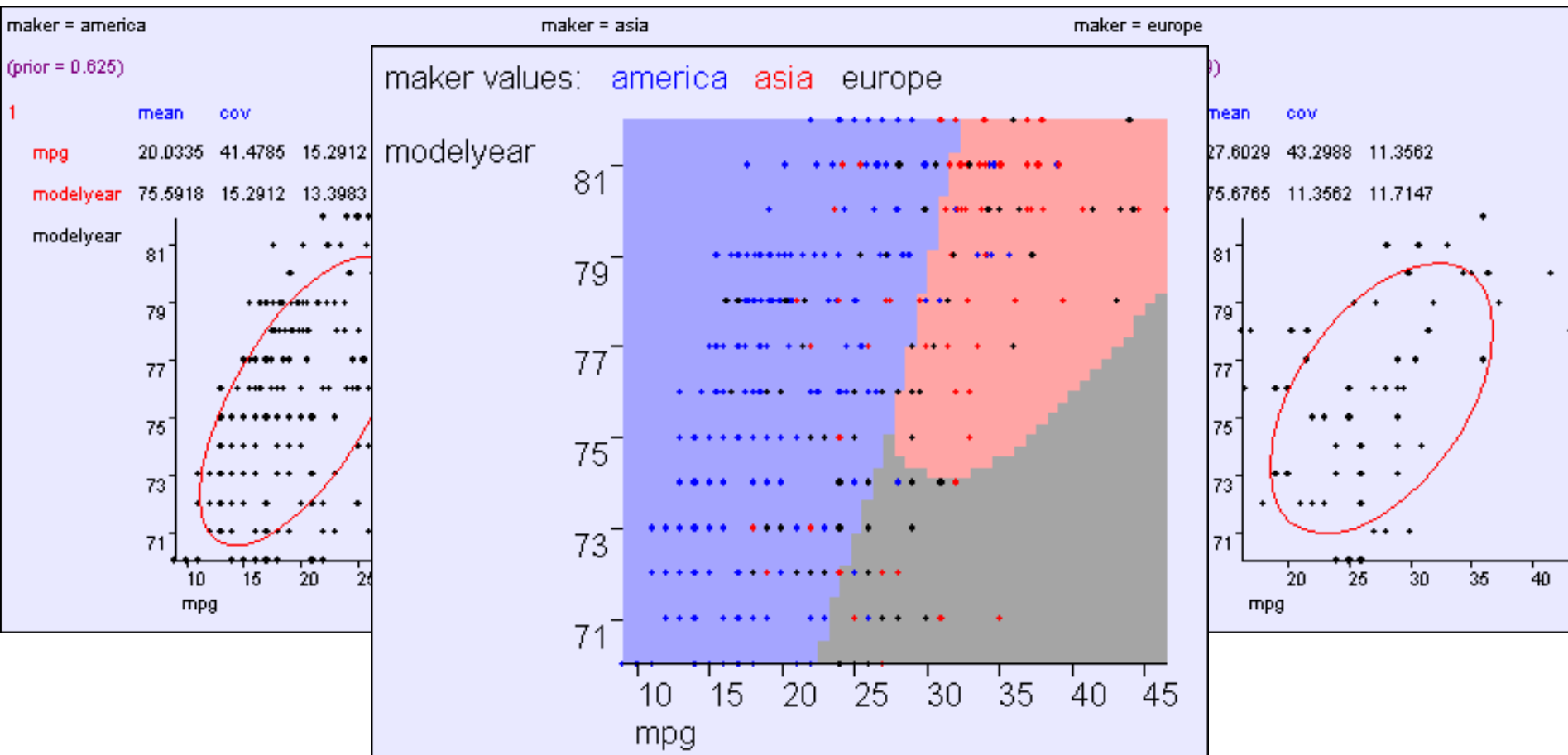
# Learning modelyear, mpg ---> maker

$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma_m^2 \end{pmatrix}$$



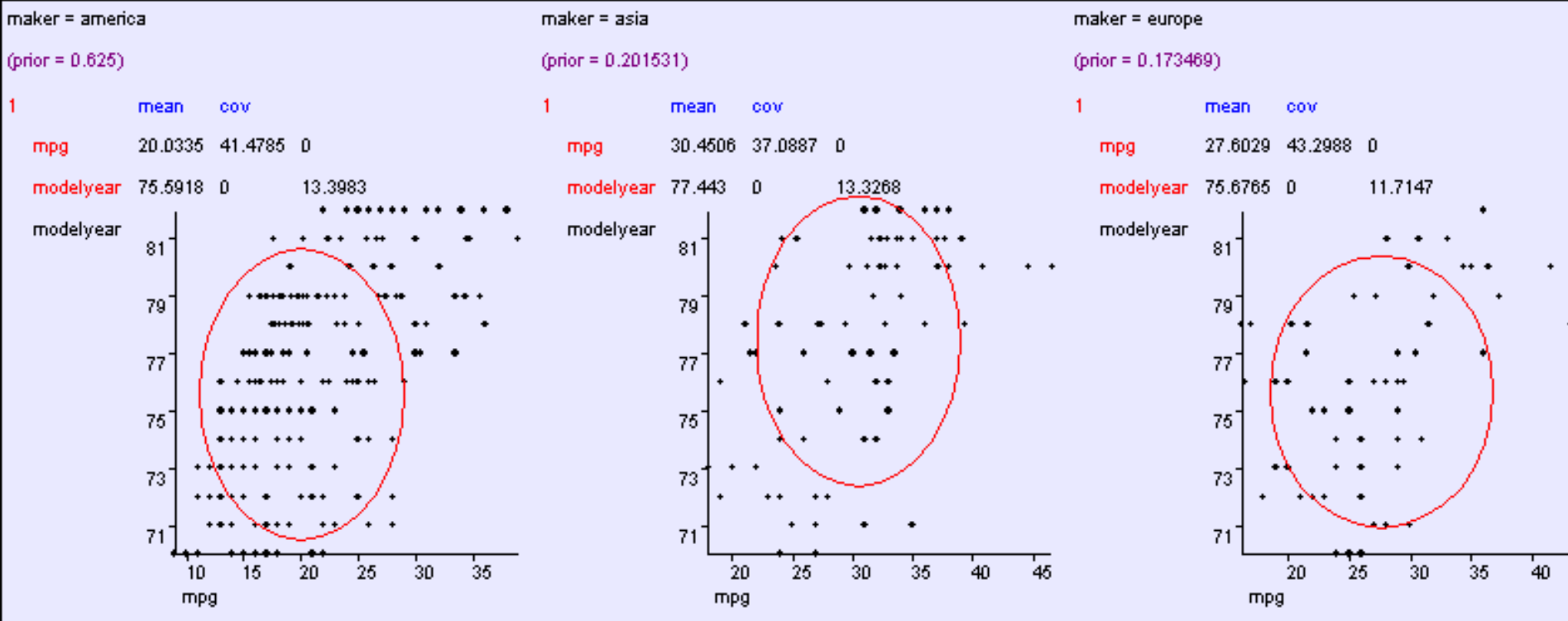
# General: $O(m^2)$ parameters

$$\Sigma = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma_m^2 \end{pmatrix}$$



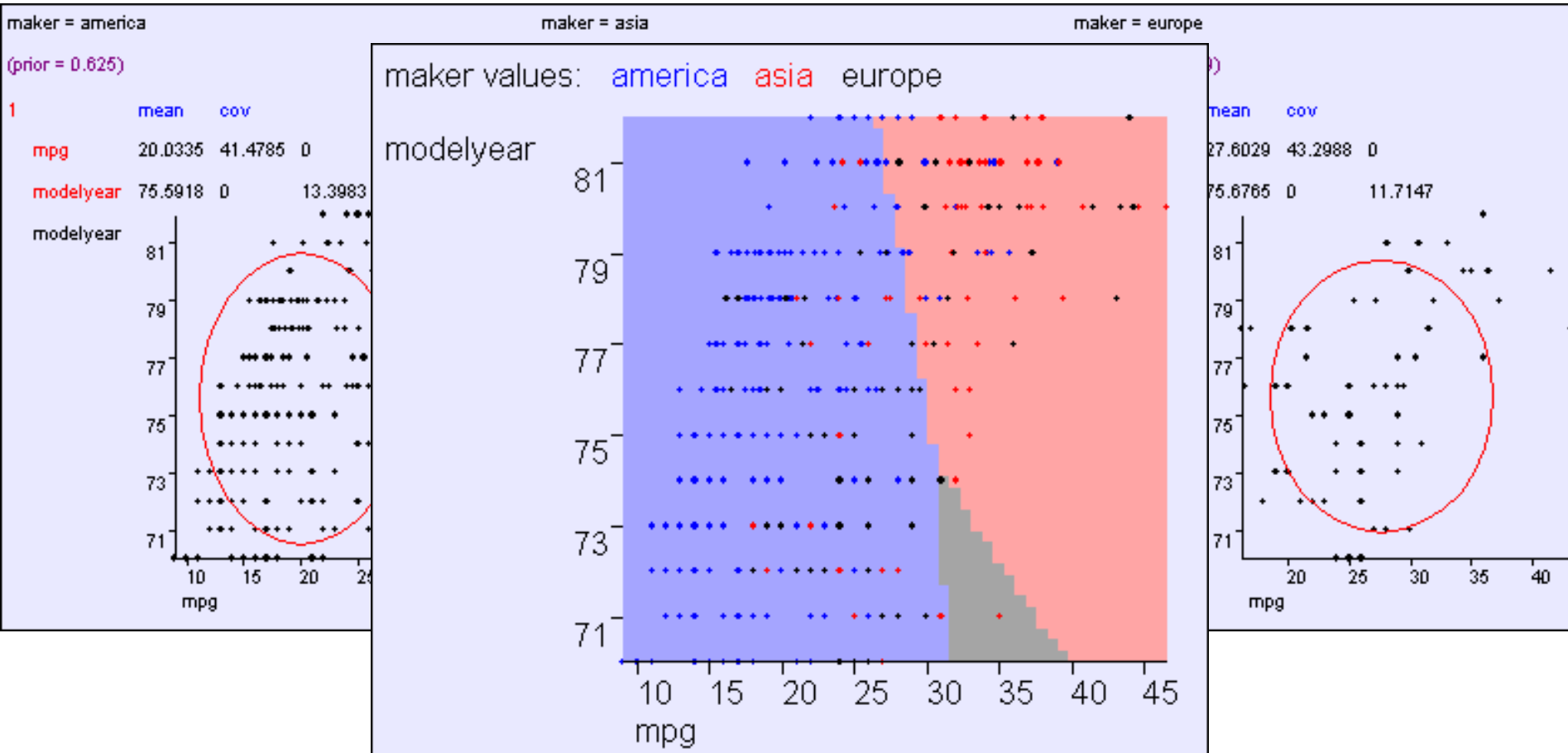
# Aligned: $O(m)$ parameters

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



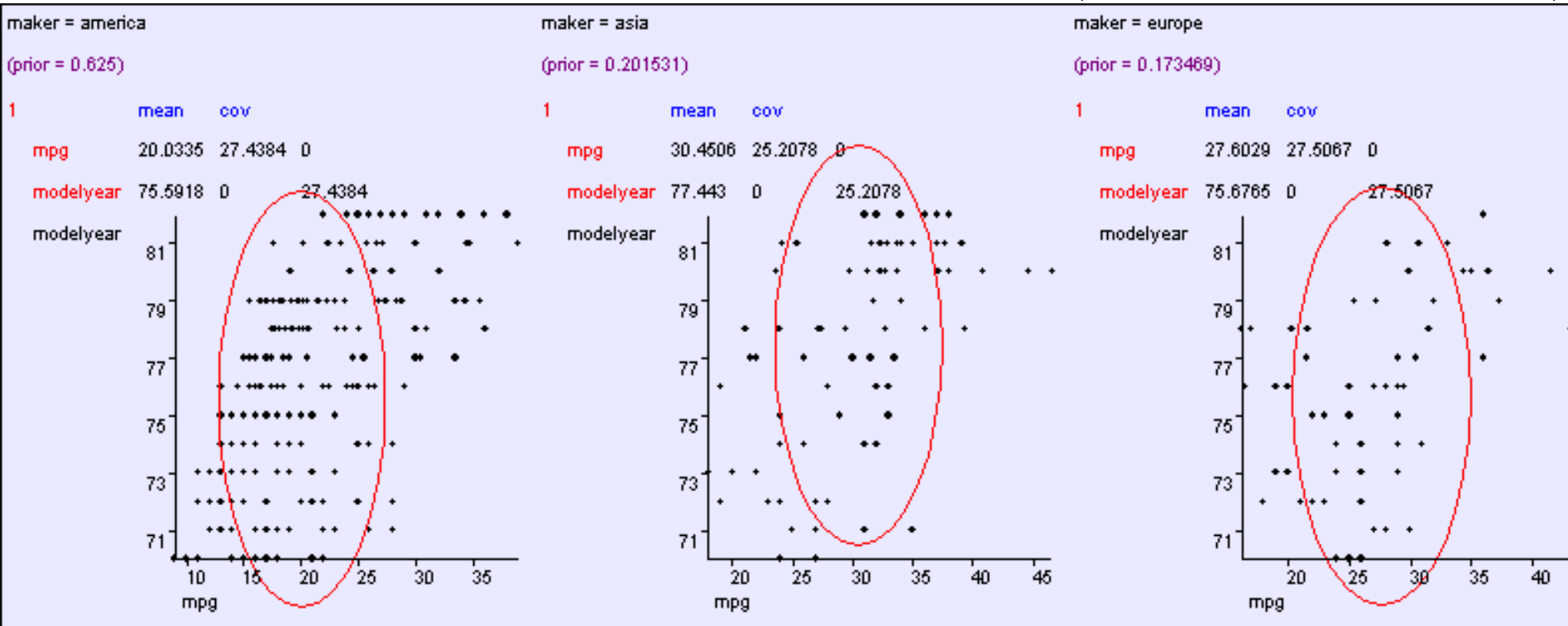
# Aligned: $O(m)$ parameters

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_2^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_3^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{m-1}^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_m^2 \end{pmatrix}$$



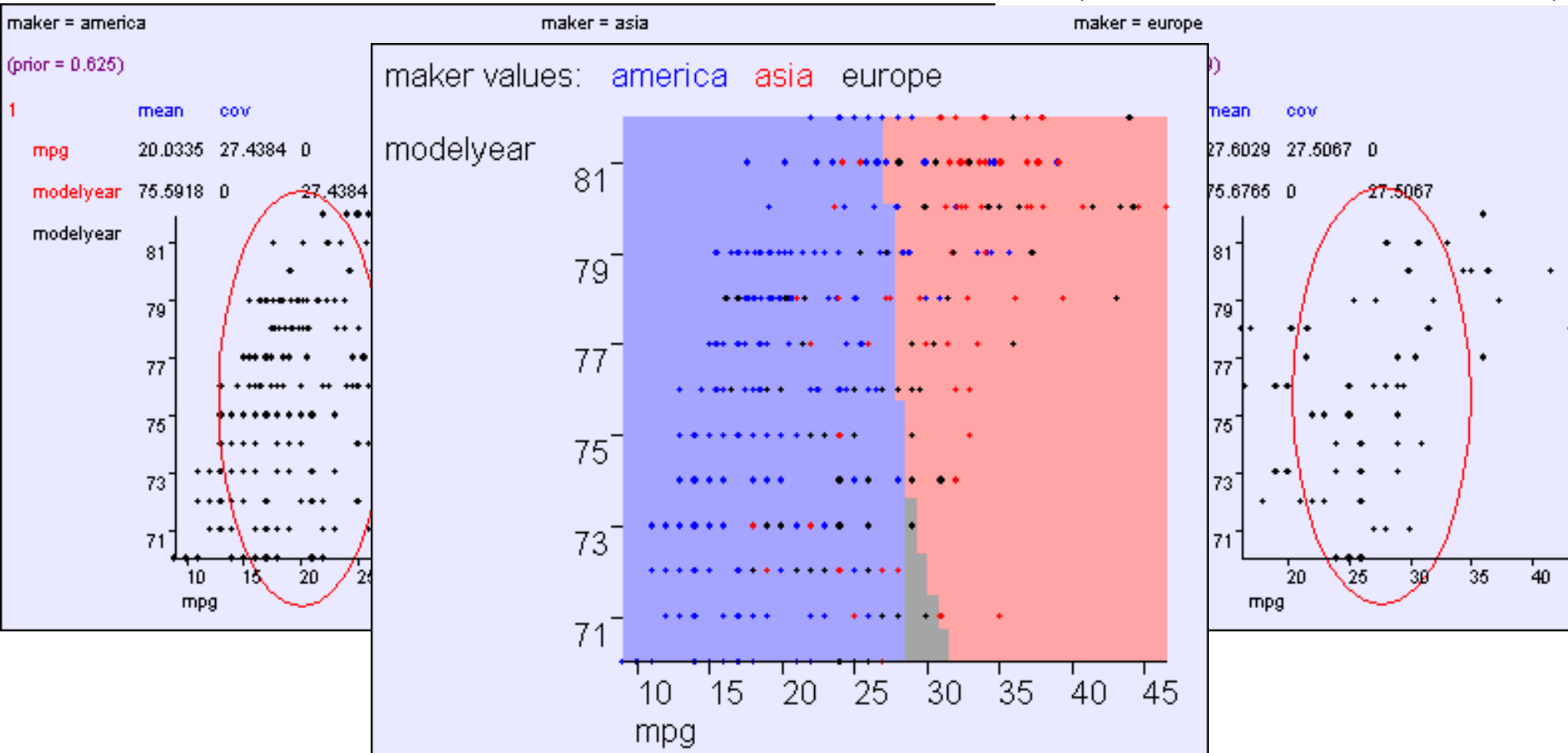
# Spherical: $O(1)$ cov parameters

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



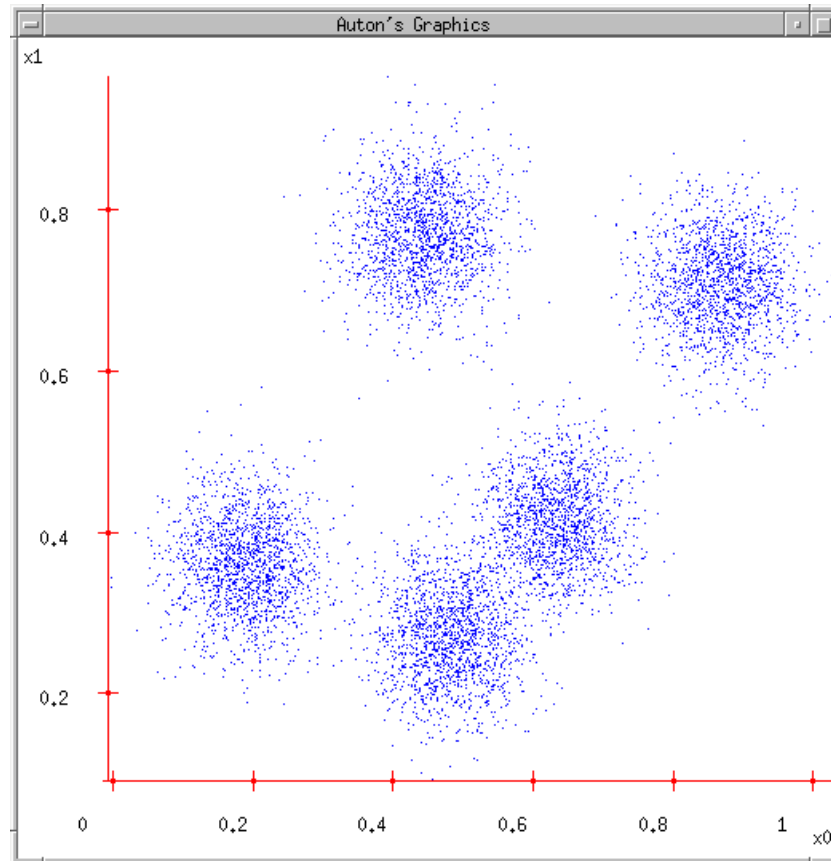
# Spherical: $O(1)$ cov parameters

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



# Density estimation

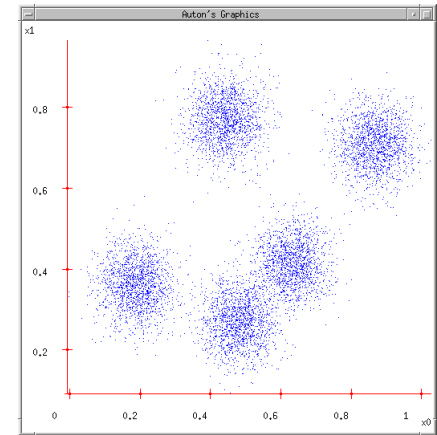
What if we want to do density estimation with multimodal or clumpy data?





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

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



# Special case: spherical Gaussians and hard assignments

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

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

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

- If each  $\mathbf{x}_j$  belongs to one class  $C(j)$  (hard assignment), marginal likelihood:

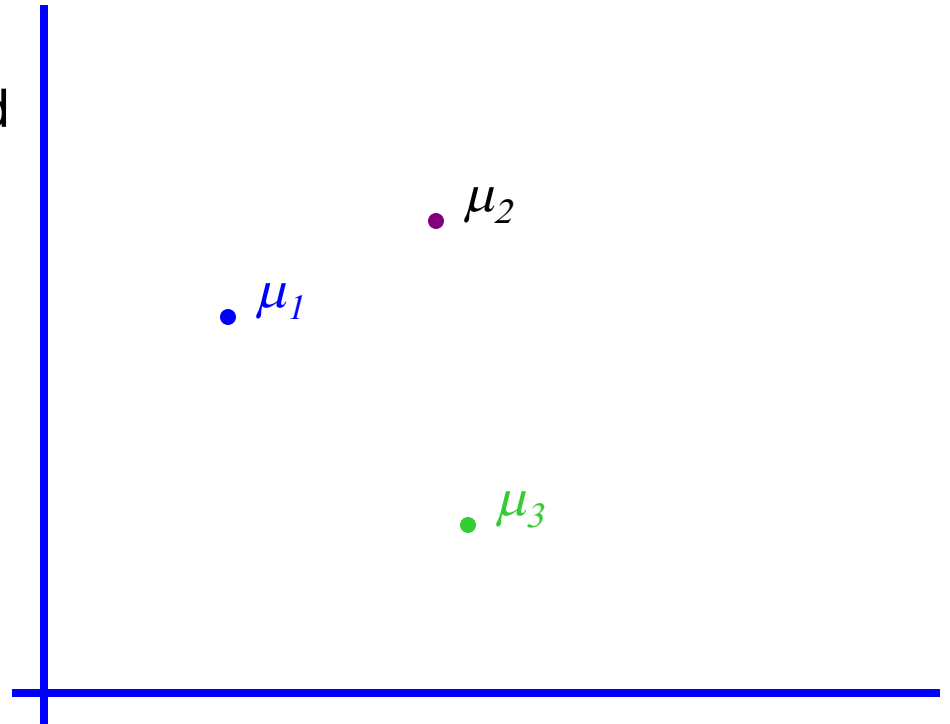
$$\prod_{j=1}^m \sum_{i=1}^k P(\mathbf{x}_j, y = i) \propto \prod_{j=1}^m \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{x}_j - \mu_{C(j)}\|^2\right]$$

- Same as K-means!!!



# The GMM assumption

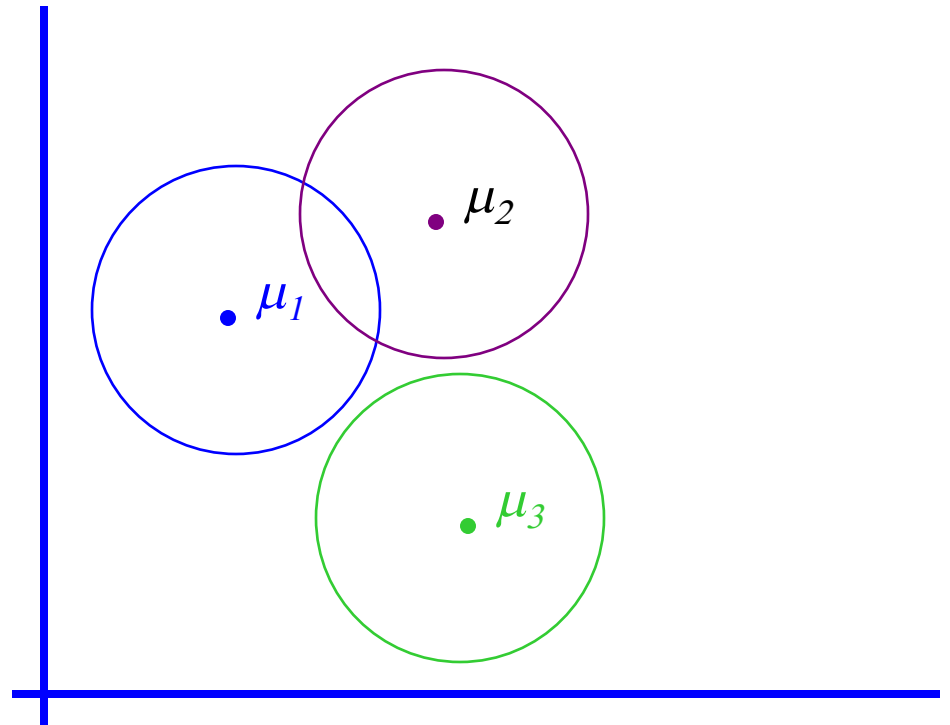
- There are  $k$  components
- Component  $i$  has an associated mean vector  $\mu_i$



# 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:

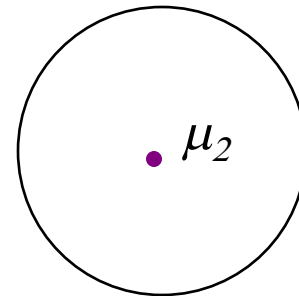


# The GMM assumption

- 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  $P(y=i)$

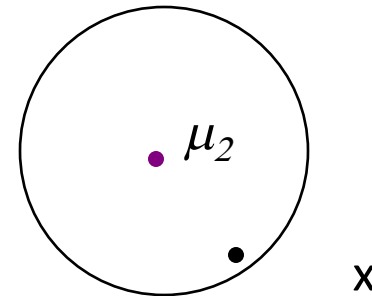


# The GMM assumption

- 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. Data point  $\sim N(\mu_i, \sigma^2 I)$

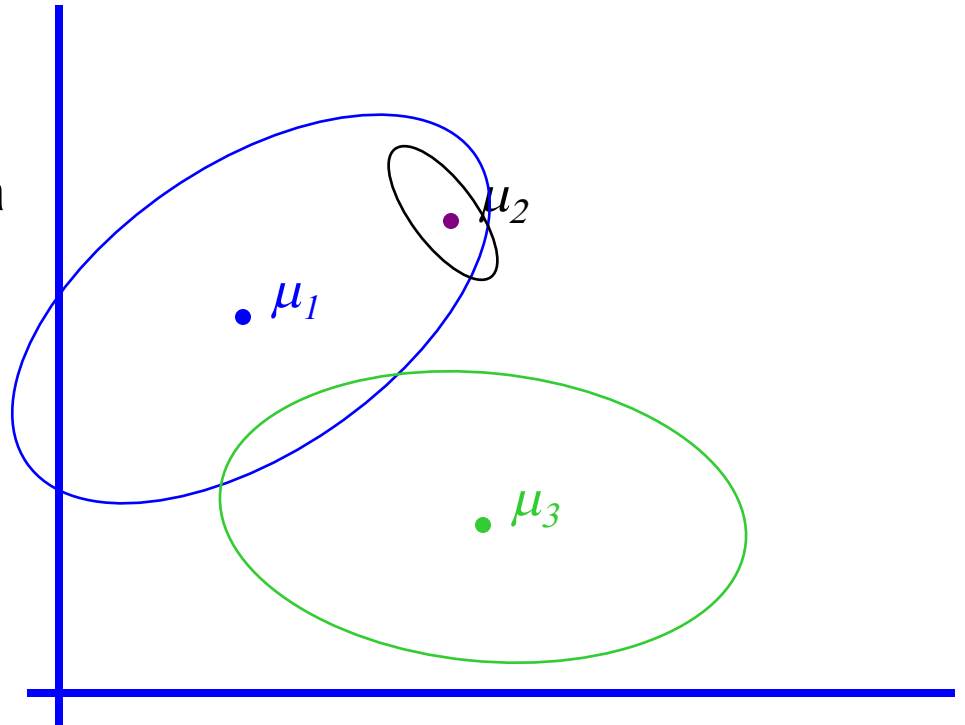


# The **General** GMM assumption

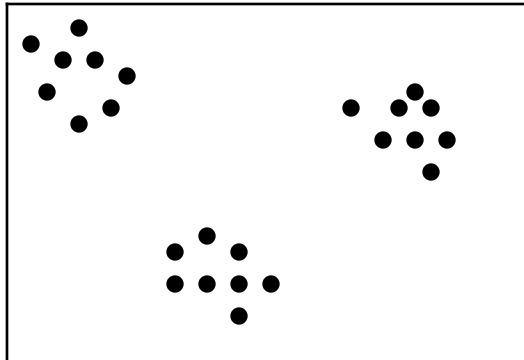
- There are  $k$  components
- Component  $i$  has an associated mean vector  $\mu_i$
- Each component generates data from a Gaussian with mean  $\mu_i$  and covariance matrix  $\Sigma_i$

Each data point is generated according to the following recipe:

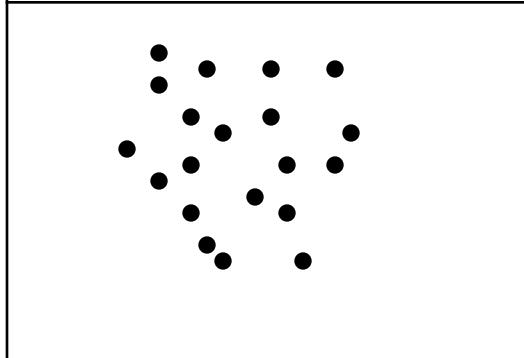
1. Pick a component at random:  
Choose component  $i$  with probability  $P(y=i)$
2. Data point  $\sim N(\mu_i, \Sigma_i)$



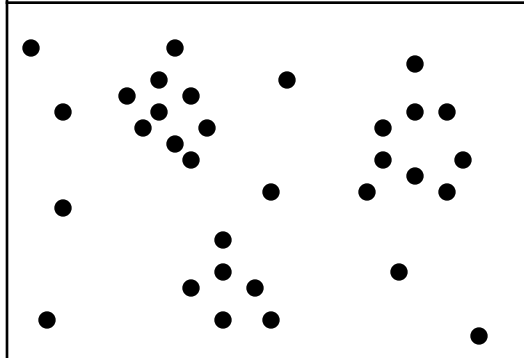
# Unsupervised Learning: not as hard as it looks



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*



# Marginal likelihood for general case

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

- Marginal likelihood:

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

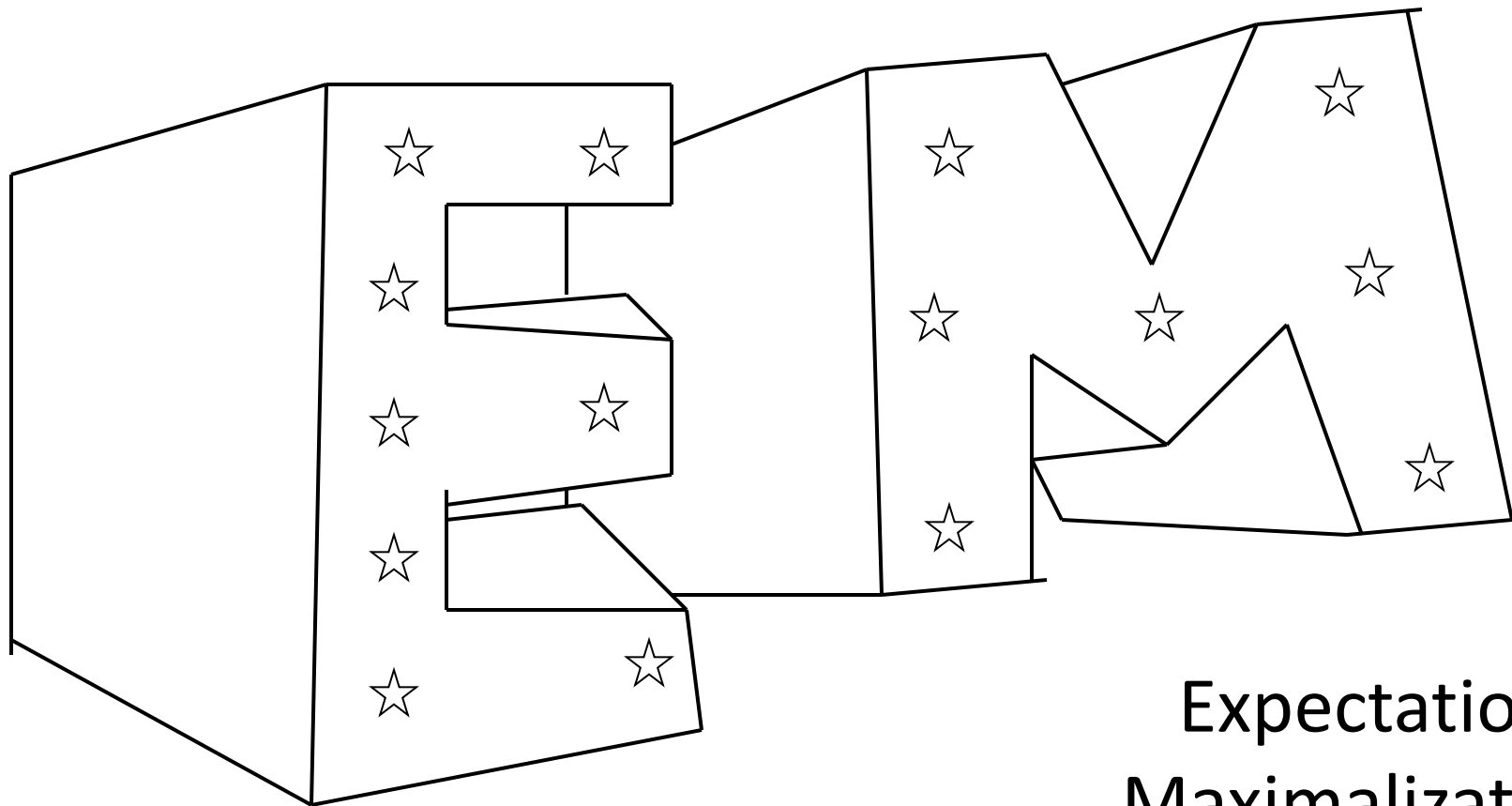
# Special case: spherical Gaussians and soft assignments

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

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

- Uncertain about class of each  $\mathbf{x}_j$  (soft assignment), marginal likelihood:

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



Expectation  
Maximalization

# 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 ?

# Trivial statistics

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

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

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

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

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

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

So if class got

A	B	C	D
14	6	9	10

$$\text{Max like } \mu = \frac{1}{10}$$

Boring, but true!

# 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

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

$$P(B) = \mu$$

$$P(C) = 2\mu$$

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

# 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 =  $d$

REMEMBER

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

$$P(B) = \mu$$

$$P(C) = 2\mu$$

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

What is the max. like estimate of  $\mu$  now?

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  $\frac{1}{2} : \mu$

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

## MAXIMIZATION

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

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

# E.M. for our trivial problem

We begin with a guess for  $\mu$

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

REMEMBER

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

$$P(B) = \mu$$

$$P(C) = 2\mu$$

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

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

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

**E-step**

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

**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.**



# E.M. Convergence

- Convergence proof based on fact that  $\text{Prob}(\text{data} \mid \mu)$  must increase or remain same between each iteration [NOT OBVIOUS]
- But it can never exceed 1 [OBVIOUS]

So it must therefore converge [OBVIOUS]

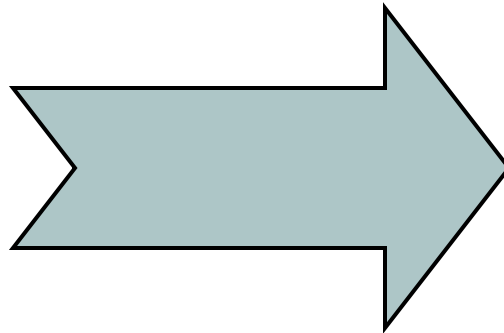
In our example,  
suppose we had

$$h = 20$$

$$c = 10$$

$$d = 10$$

$$\mu^{(0)} = 0$$



Convergence is generally linear: 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  $\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_m$

We know there are  $k$  classes

We know  $P(y_1) P(y_2) P(y_3) \dots P(y_k)$

We don't know  $\mu_1 \mu_2 \dots \mu_k$

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

$$= p(x_1 \dots x_m \mid \mu_1 \dots \mu_k)$$

$$= \prod_{j=1}^m p(x_j \mid \mu_1 \dots \mu_k) \quad \text{🗨️}$$

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

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

# EM for simple case of GMMs:

## The E-step

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

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

# EM for simple case of GMMs:

## The M-step

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

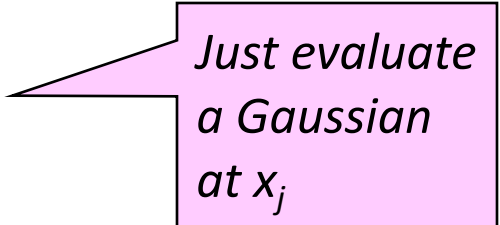
$$\mu_i = \frac{\sum_{j=1}^m P(y=i | x_j) x_j}{\sum_{j=1}^m P(y=i | x_j)}$$

# E.M. for GMMs

## E-step

Compute “expected” classes of all data points for each class

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



*Just evaluate  
a Gaussian  
at  $x_j$*

## M-step

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

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