

Unsupervised learning or Clustering (cont.) – *EM algorithm* K-means Gaussian mixture models

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

April 5th, 2006

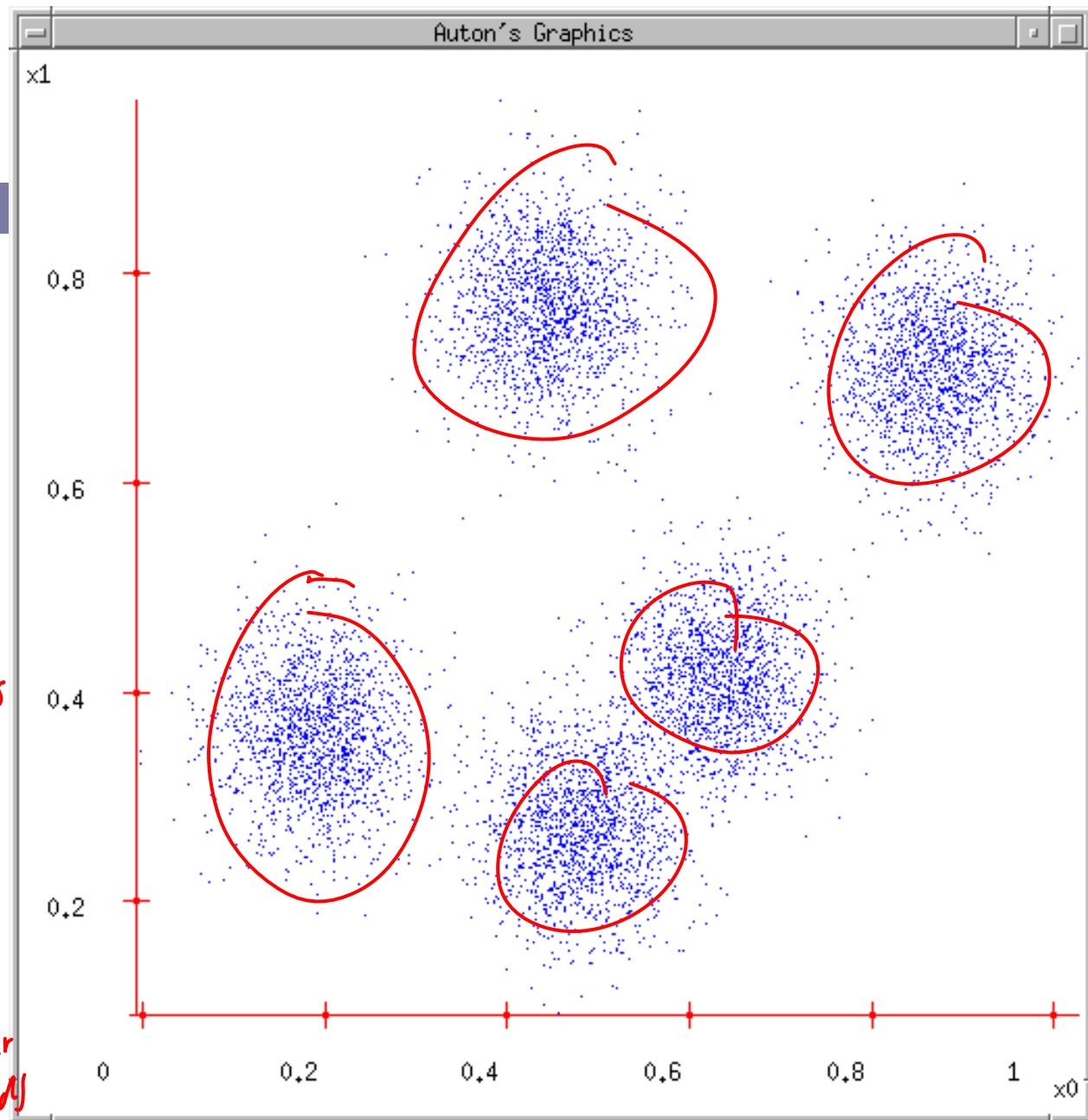
Some Data

nobody tells
you what the
clusters are...

classification
give $x \rightarrow$ must tell you class

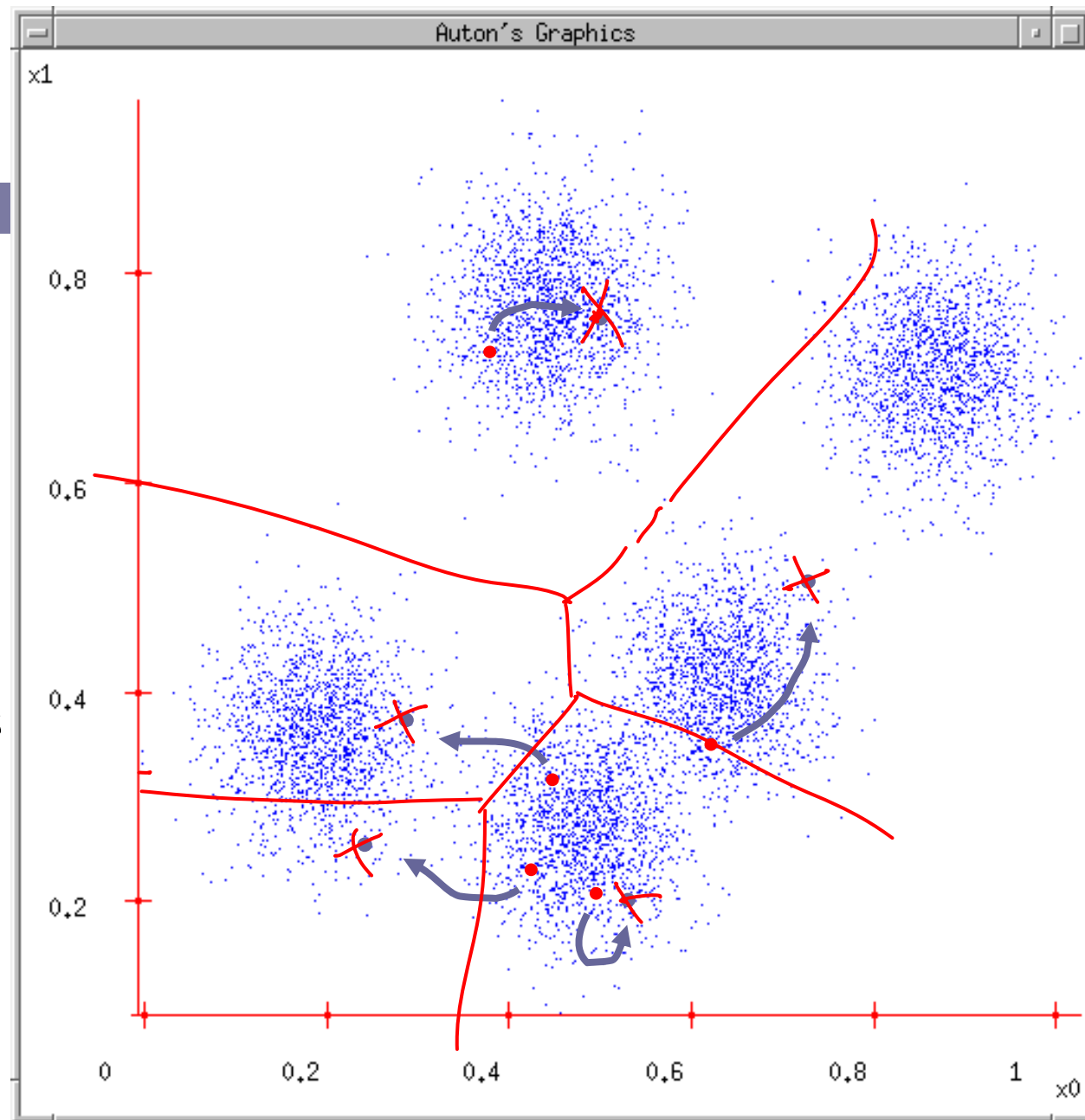
train data (supervised)
 (x_i, y_i) know label

cluster
give $x \rightarrow$ must tell you cluster
train data (unsupervised)
 (x_i) no labels



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!



K-means

- Randomly initialize k centers

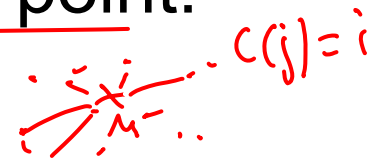
- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$ *iteration*

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center: *centers for iteration t*

- $C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$
point j (i)

- **Recenter:** μ_i becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2$



- Equivalent to $\mu_i \leftarrow$ average of its points!

What is K-means optimizing?

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

$$\square F(\mu, C) = \sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2 \quad \text{loss function}$$
$$= \sum_{i=1}^K \sum_{j: C(j)=i} \|\mu_i - x_j\|^2$$

centers (pointing to μ)
allocation (pointing to C)

- Optimal K-means:

$$\square \min_{\mu} \min_C F(\mu, C)$$

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 μ , optimize C

fix $\mu = \bar{\mu}$

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

$$C(j) \leftarrow \operatorname{argmin}_i \| \bar{\mu}_i - x_j \|^2$$

achieve minimum..

exactly the K-classification step in K-means

$$F(\mu^{(t+1)}, C^{(t+1)})$$

fix $\mu^{(t)}$
opt $C^{(t+1)}$

$$F(\mu^{(t)}, C^{(t+1)}) \leq F(\mu^{(t)}, C^{(t)})$$

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 $C^{(t+1)}$
 opt to get $\mu^{(t+1)}$
 $F(C^{(t+1)}, \mu^{(t+1)}) \leq F(C^{(t+1)}, \mu^{(t)})$

- Fix C, optimize μ

$C \leftarrow \bar{C}$

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

μ_i should be centroid

recentering step in K-means

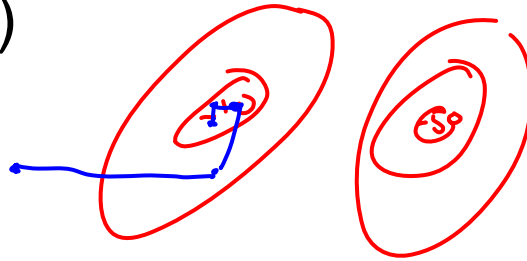
Coordinate descent algorithms

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

Want: $\min_a \min_b F(a,b)$

Coordinate descent:

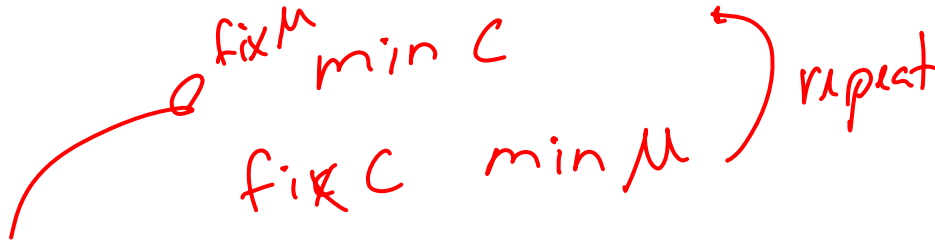
- fix a, minimize b
- fix b, minimize a
- repeat



Converges!!!

$$\exists K: \forall a, b \quad f(a,b) \geq K$$

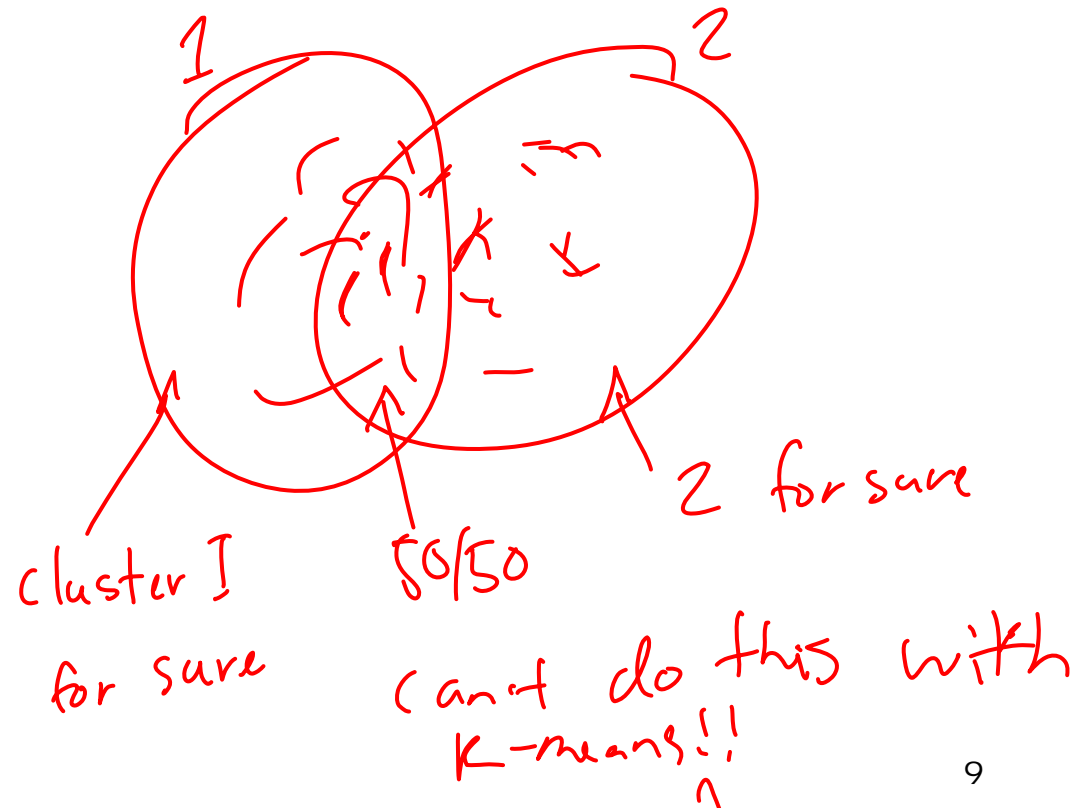
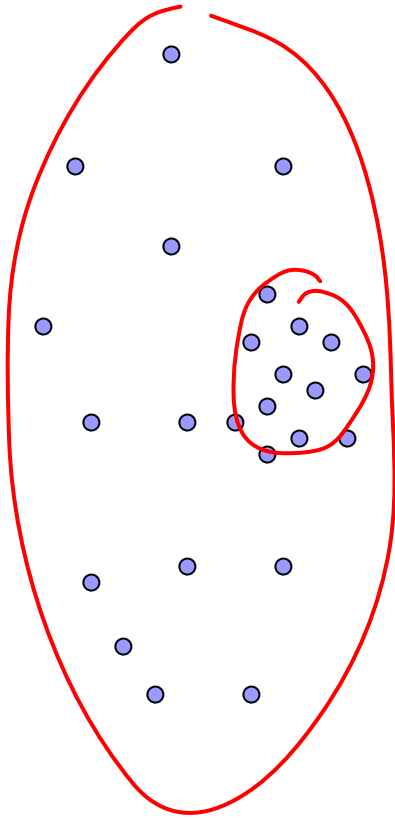
- if F is bounded
- to a (often good) local optimum
 - as we saw in applet (play with it!)



K-means is a coordinate descent algorithm!

(One) bad case for k-means

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

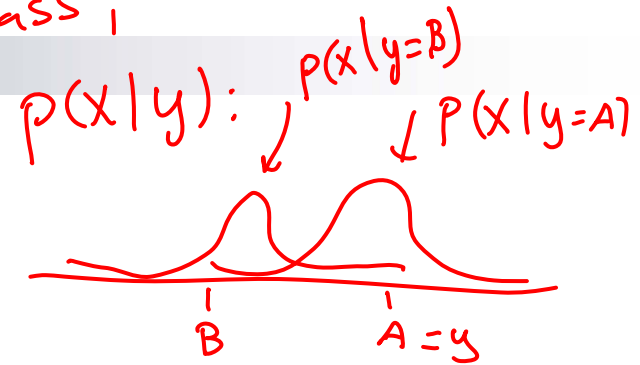


Gaussian Bayes Classifier

Reminder

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

class $y=i$ *input* \mathbf{x}_j



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

prop. $\frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}}$ *gaussian for class i* *prior* $P(y = i)$

$$p(x|y=i) \sim N(\mu_i, \Sigma_i)$$

$$\mu_i = \begin{pmatrix} 10 \\ -5 \\ 200 \\ \vdots \end{pmatrix}$$

→ give me $x_j = 98$

$$\left. \begin{array}{l} P(y=A|98) \propto 0.4 \\ P(y=B|98) \propto 0.1 \end{array} \right\} \begin{array}{l} 0.8 \\ 0.2 \end{array}$$

renormalize

$$\Sigma_i = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ & \sigma_2^2 & \sigma_{23} \\ & & \sigma_3^2 \end{pmatrix}$$

Var. $\sigma_1^2, \sigma_2^2, \sigma_3^2$ *Co-var.* $\sigma_{12}, \sigma_{13}, \sigma_{23}$

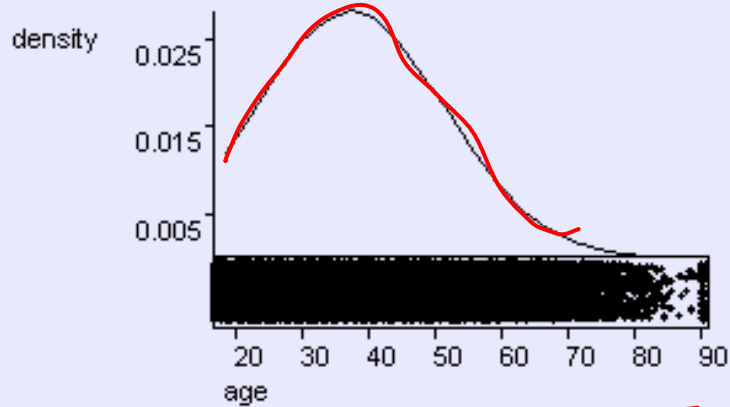
Predicting wealth from age

wealth = poor

(prior = 0.760718)

1 mean cov

age 37.374 198.935

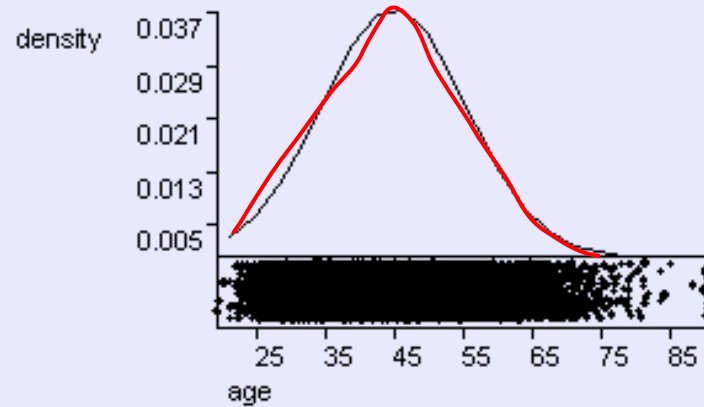


wealth = rich

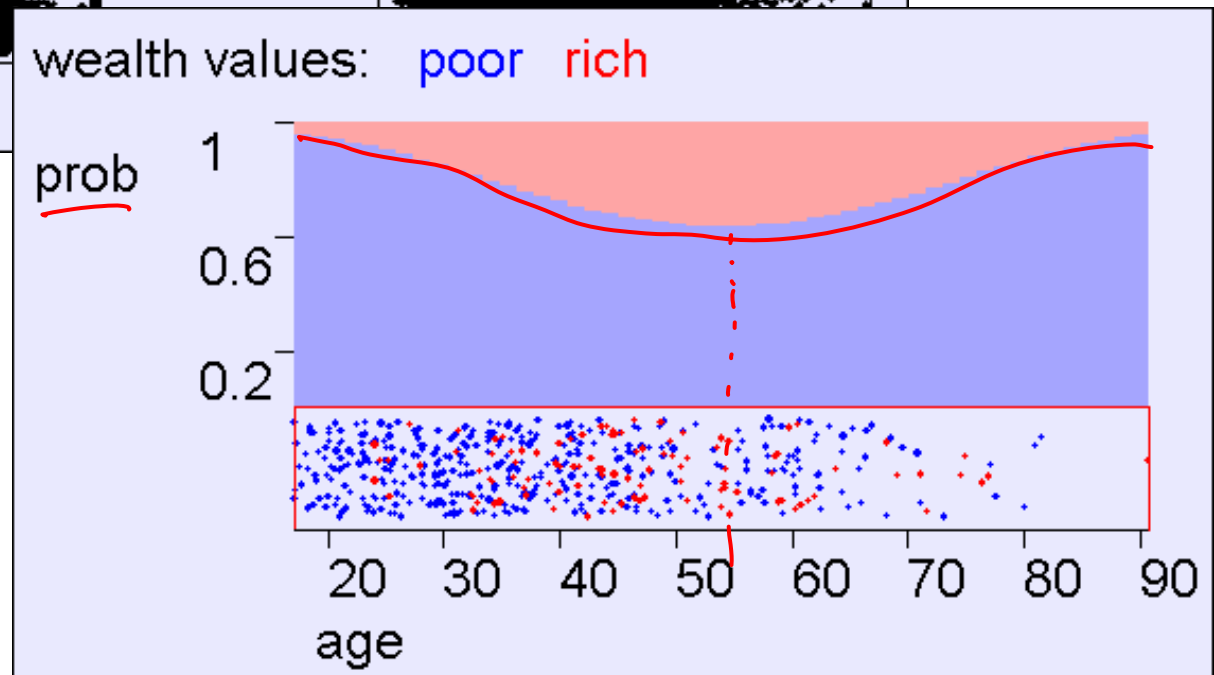
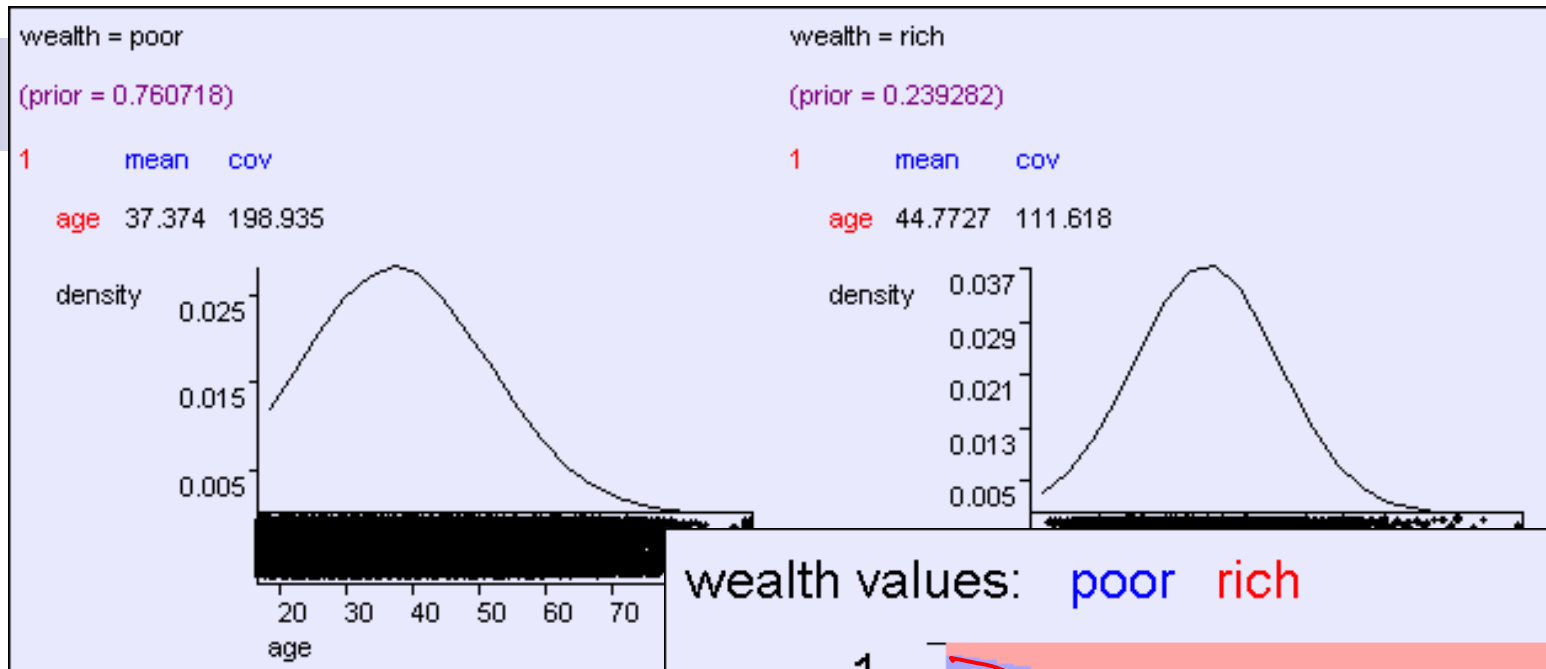
(prior = 0.239282)

1 mean cov

age 44.7727 111.618



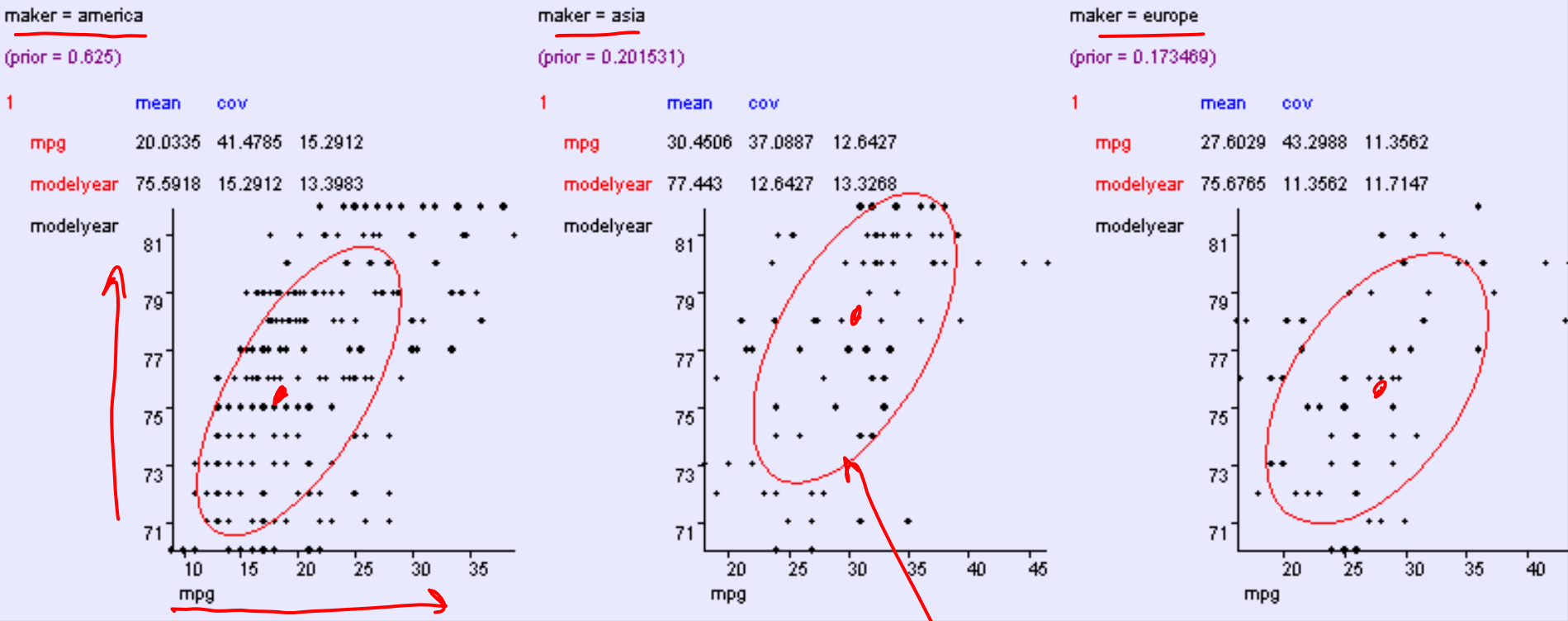
Predicting wealth from age



Learning modelyear, mpg ---> maker

Σ is 2x2 [2 continuous features]

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



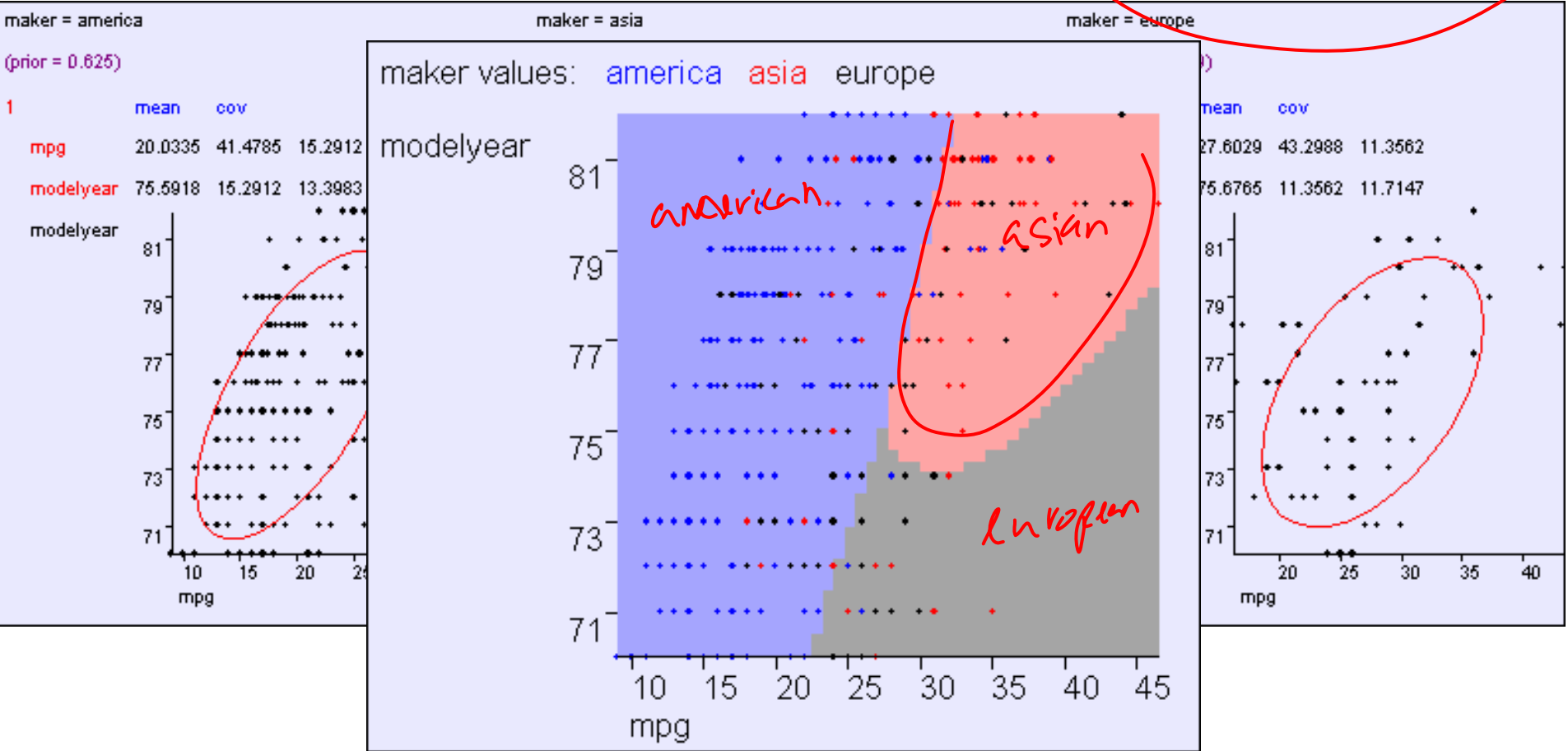
tilted bc cause $\sigma_{ij} \neq 0$

General: $O(m^2)$

parameters

m features

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

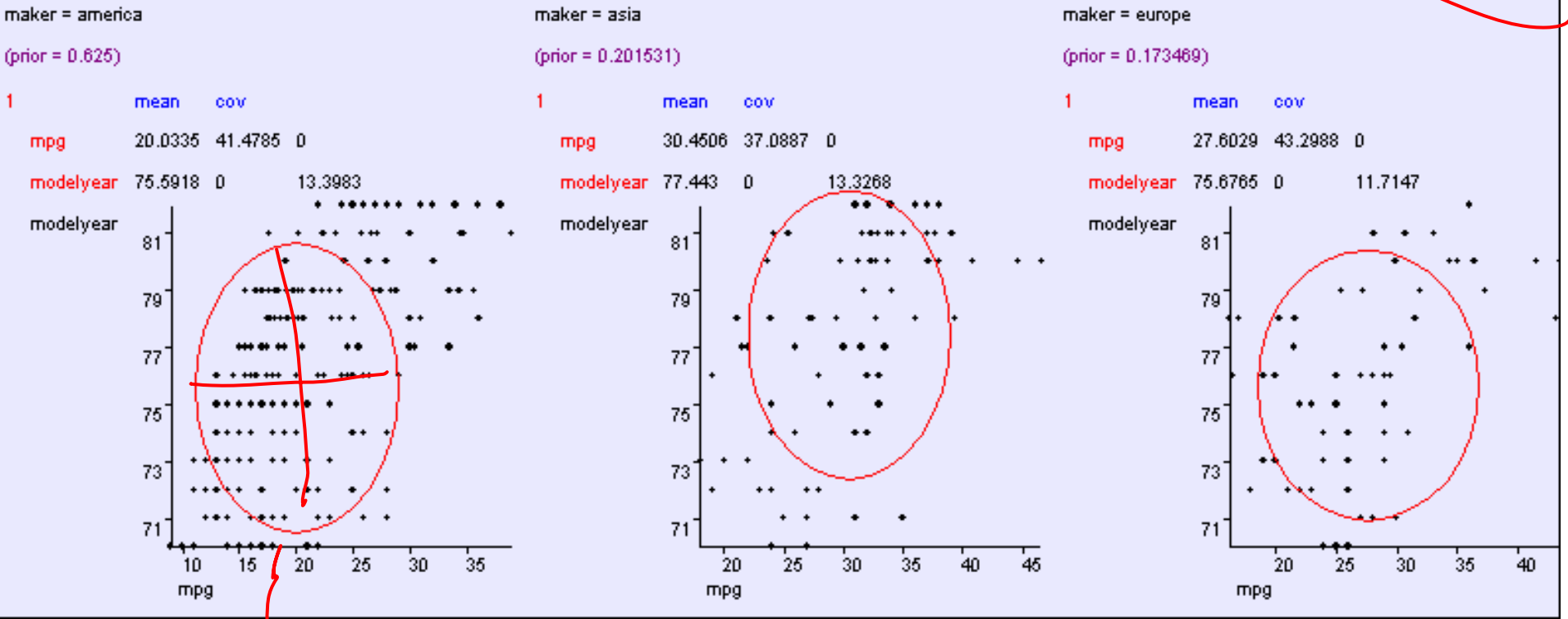


Aligned: $O(m)$ parameters

only variances

one per class

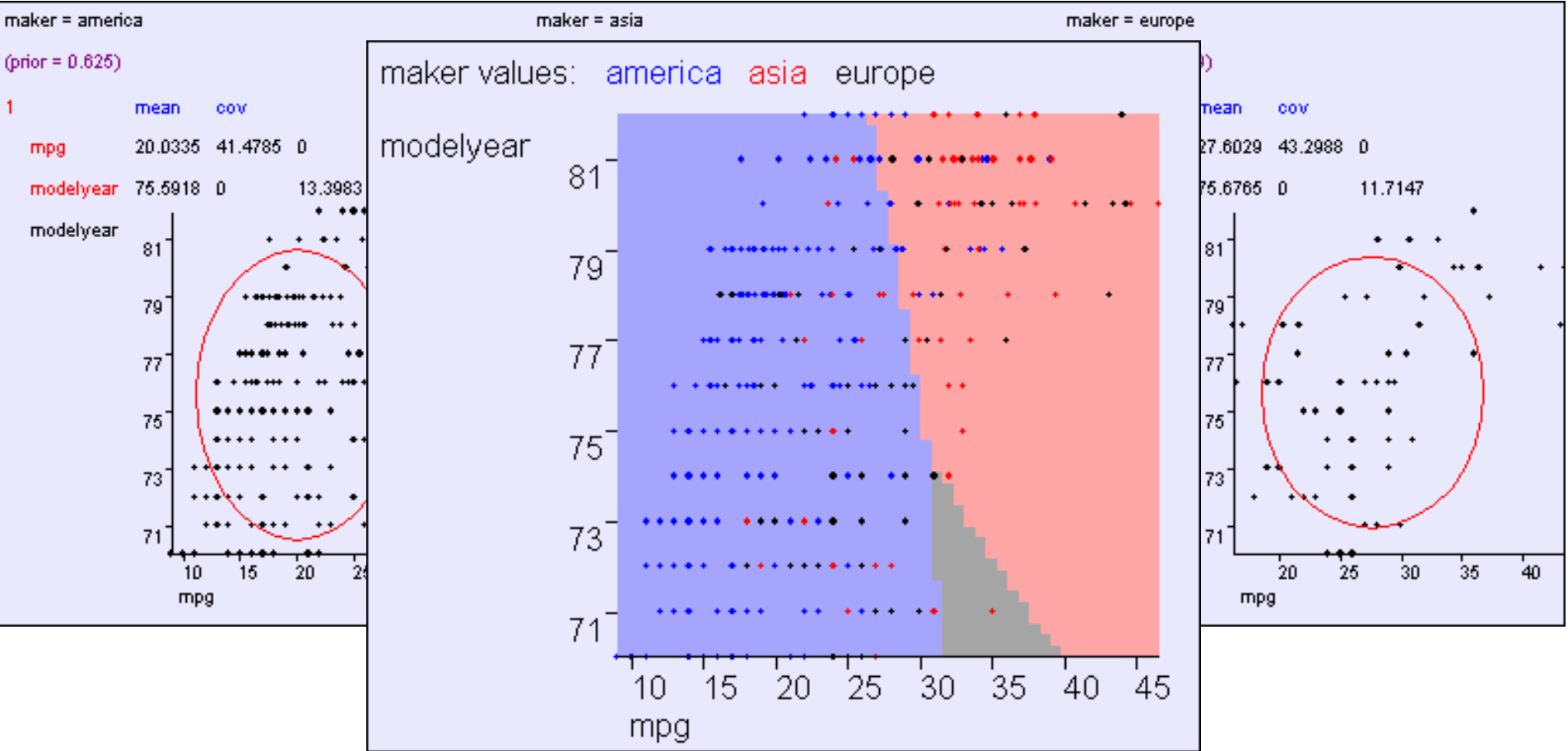
$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma^2_m \end{pmatrix}$$



axis aligned

Aligned: $O(m)$ parameters

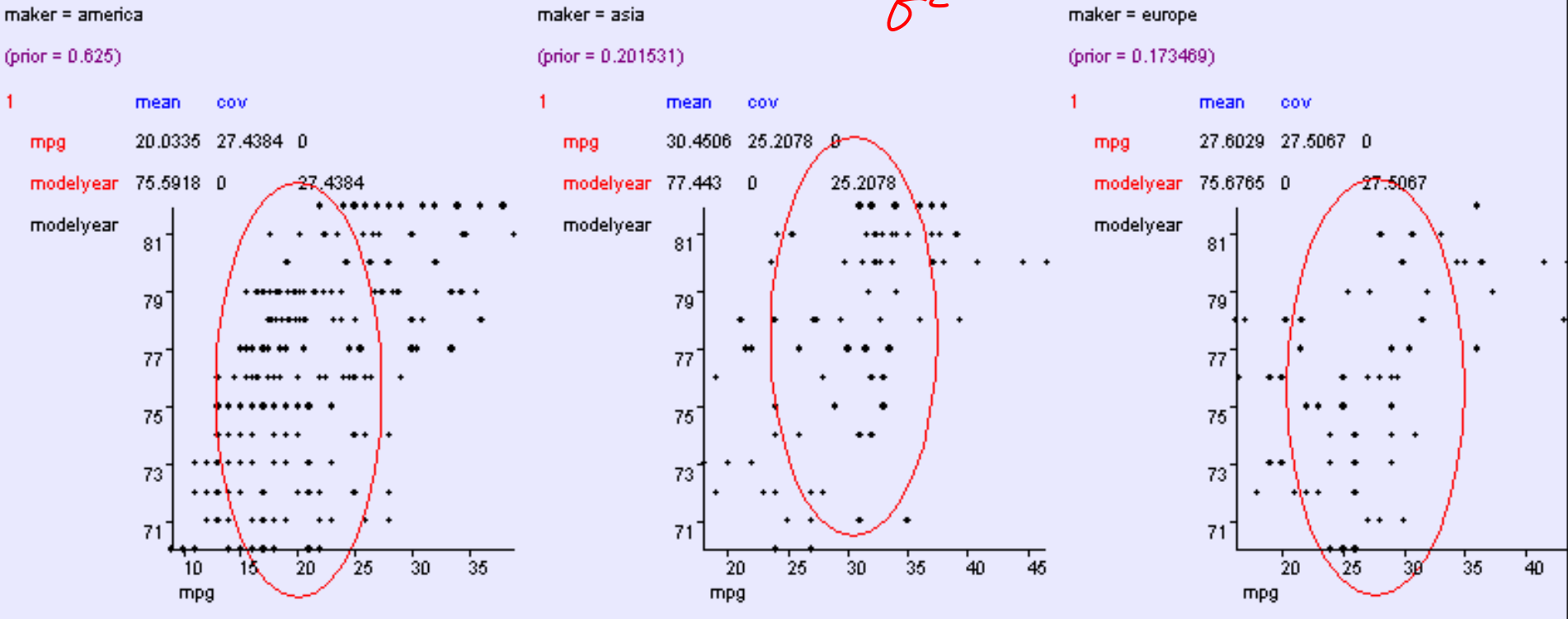
$$\Sigma = \begin{pmatrix} \sigma^2_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma^2_2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma^2_3 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2_{m-1} & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma^2_m \end{pmatrix}$$



Spherical: $O(1)$ cov parameters

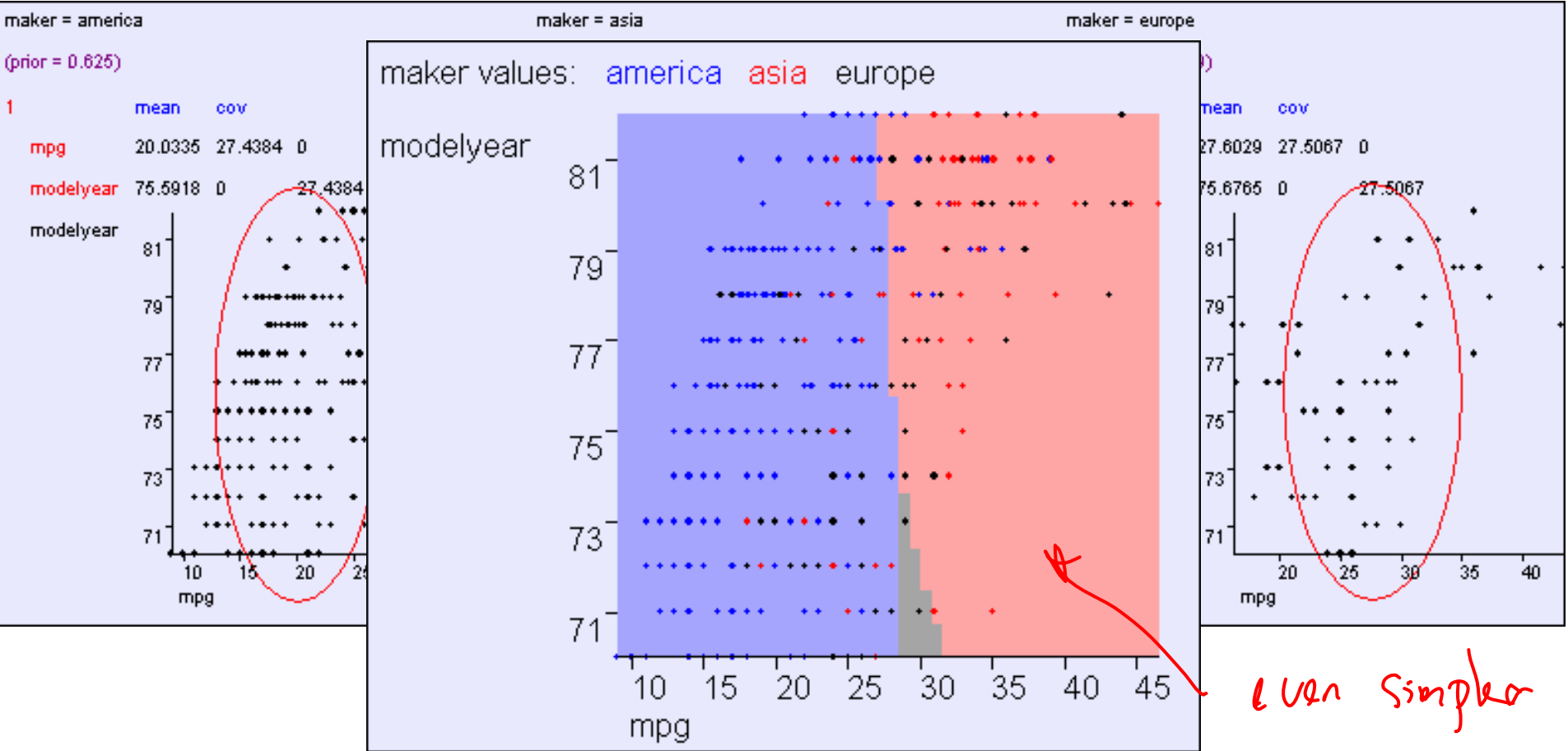
$$\Sigma = \begin{pmatrix} \sigma^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma^2 \end{pmatrix}$$

all features same variance σ^2



Spherical: $O(1)$ cov parameters

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



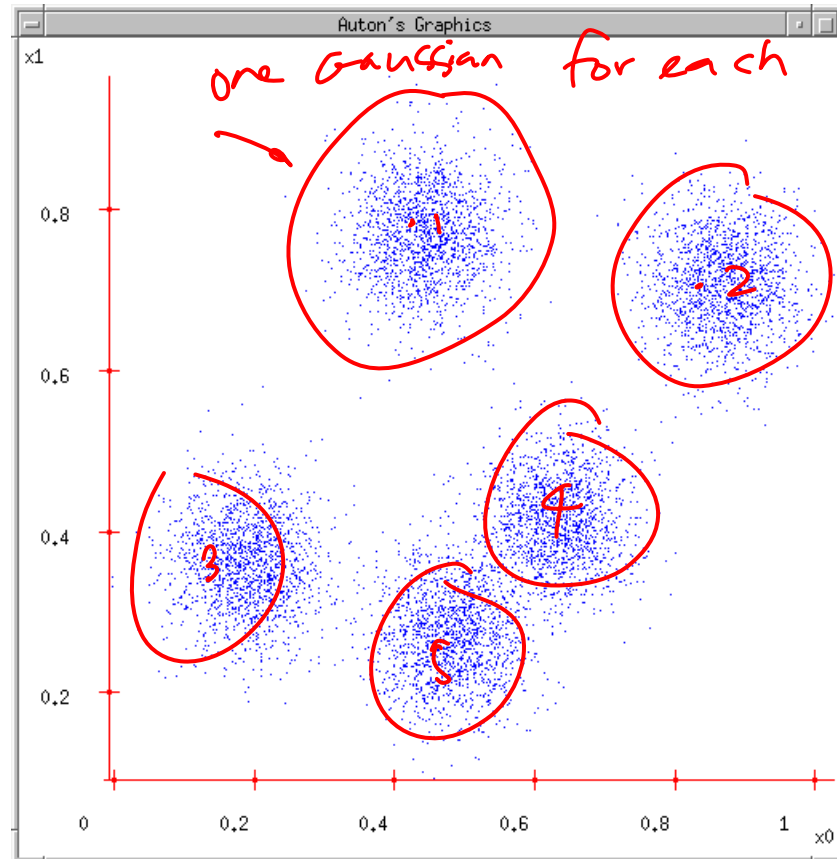
Next... back to Density Estimation

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

$$X \sim \sum_{i=1}^K w_i N(\mu_i, \Sigma_i)$$

$y \leftarrow$ don't know
point has $x, y = i$

$$X \sim N(\mu_i, \Sigma_i)$$

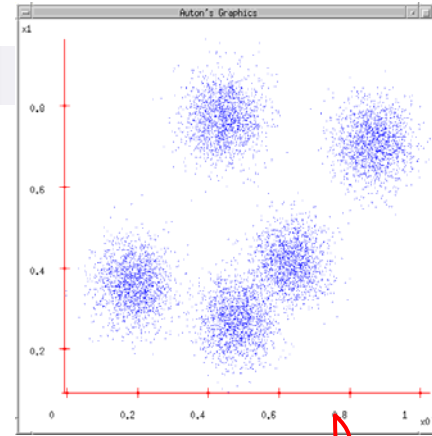


one Gaussian for each bump.

But we don't see class labels!!!

- MLE: (x_j, y_j)

$$\begin{aligned} \arg \max \prod_j P(y_j, x_j) \\ = \arg \max \log \prod_j P(y_j, x_j) &= \arg \max \sum_j \log P(y_j, x_j) \end{aligned}$$



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

almost always nice !!
 only observe x_j
 closed form or convex

- Maximize marginal likelihood:

$$\begin{aligned} \arg \max \prod_j P(x_j) &= \arg \max \prod_j \sum_{i=1}^k P(y_j=i, x_j) \\ &= \arg \max \log \prod_j P(x_j) \\ &= \arg \max \sum_j \log \sum_{i=1}^k P(y_j=i, x_j) \end{aligned}$$

log sum is almost never nice !!
 i20

Special case: spherical Gaussians and hard assignments

$$P(\mathbf{x}_j | y = i) = \frac{1}{(2\pi)^{m/2} \|\boldsymbol{\Sigma}_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x}_j - \boldsymbol{\mu}_i)\right] \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

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

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

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

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

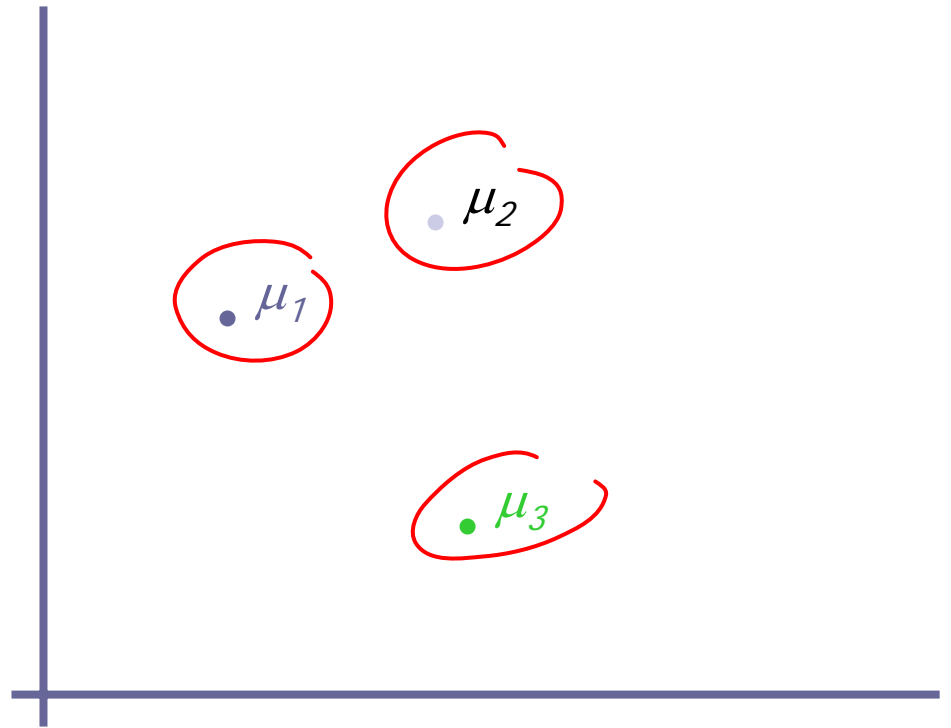
$$= \sum_j \|\mathbf{x}_j - \boldsymbol{\mu}_{C(j)}\|^2 + O(1)$$

$\Rightarrow \forall j \ P(y=i | \mathbf{x}_j)$
 $= \begin{cases} 1, & \text{otherwise} \\ 0, & C(j) \neq i \end{cases}$

- Same as K-means!!!

The GMM assumption

- There are k components
- Component i has an associated mean vector μ_i

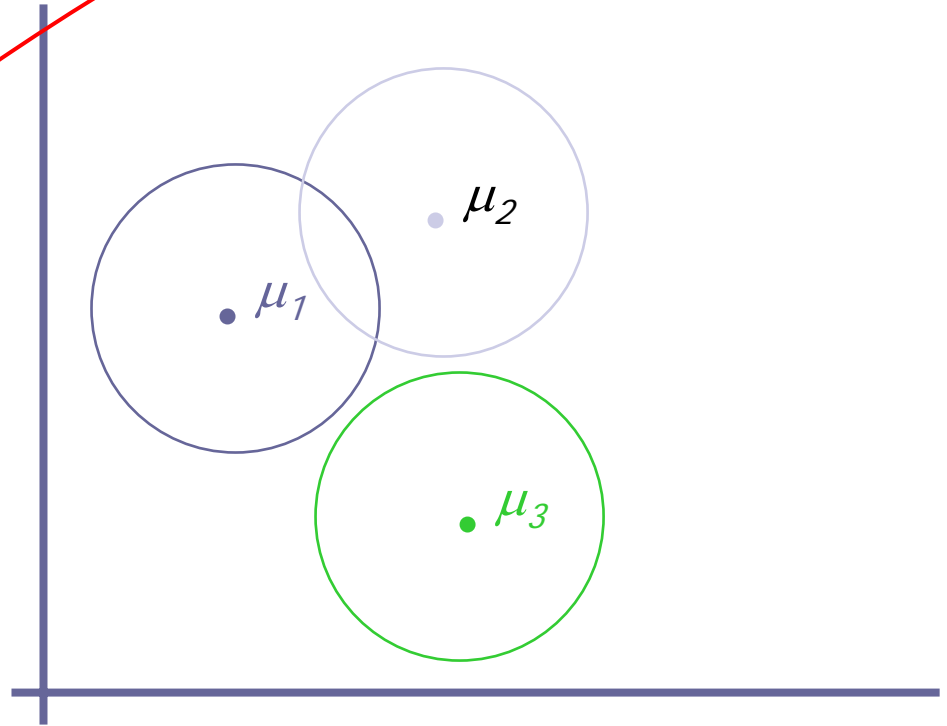


The GMM assumption

$$\begin{pmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \sigma_3^2 \end{pmatrix}$$

- There are k components
- Component i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix $\sigma^2 I$ *Spherical*

Each data point is generated according to the following recipe:

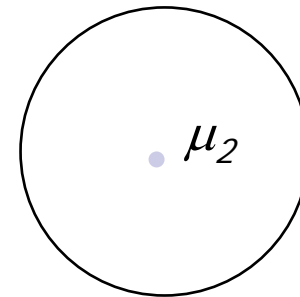


The GMM assumption

- There are k components
- Component i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_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 $P(y=i)$

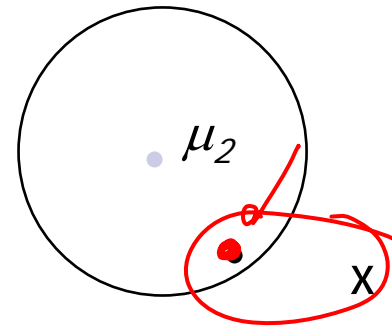


The GMM assumption

- There are k components
- Component i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_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 $\sim N(\mu_i, \sigma^2 I)$ ✓

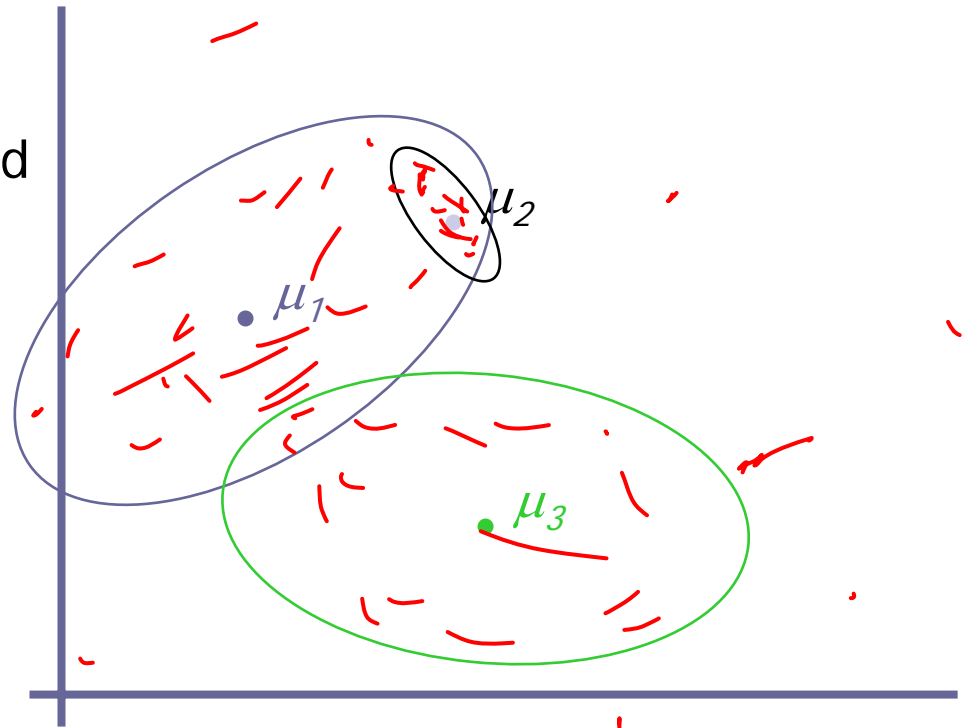


The **General** GMM assumption

- There are k components
- Component i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix Σ_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 $\sim N(\mu_i, \Sigma_i)$

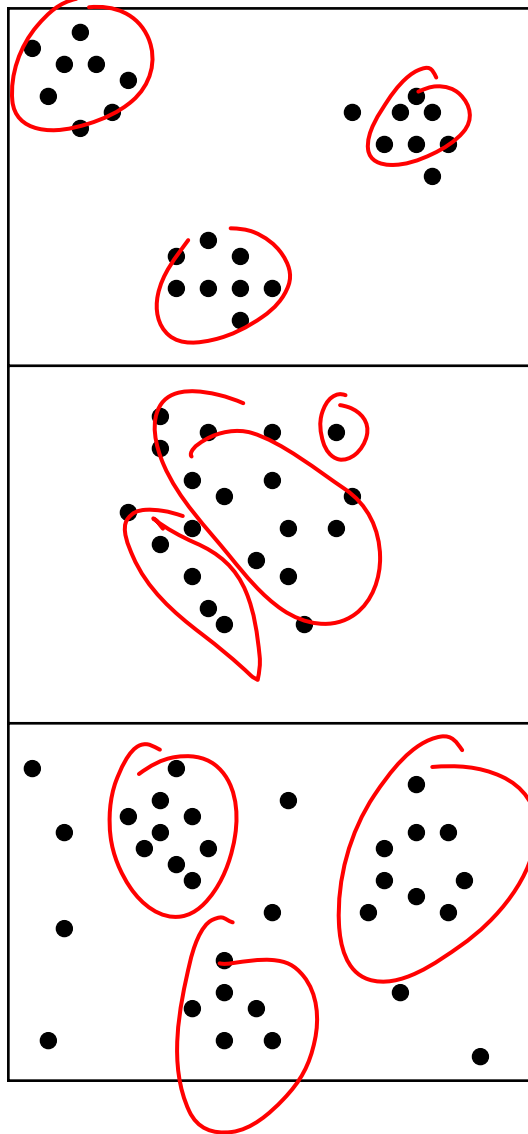


selecting cluster at random with prob. $p(y=i)$

pick location accord. gaussian

Unsupervised Learning: not as hard as it looks

given n points
re cover μ_i, Σ_i
 $i=1 \dots k$



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

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

■ Marginal likelihood:

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

$P(\mathbf{x}|y=i)$ $P(y=i)$

Special case 2: spherical Gaussians and soft assignments

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

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

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

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

log

$$\sum_{j=1}^m \log \sum_{i=1}^k e^{\left[-\frac{1}{2\sigma^2} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2\right]} P(y=i)$$

Unsupervised Learning: Mediumly Good News

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

Suppose x 's are 1-dimensional.

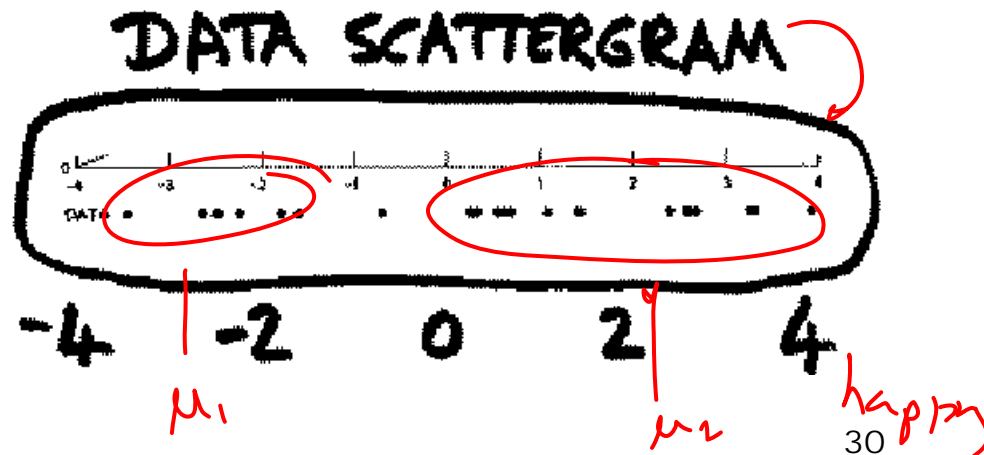
(From Duda and Hart)

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

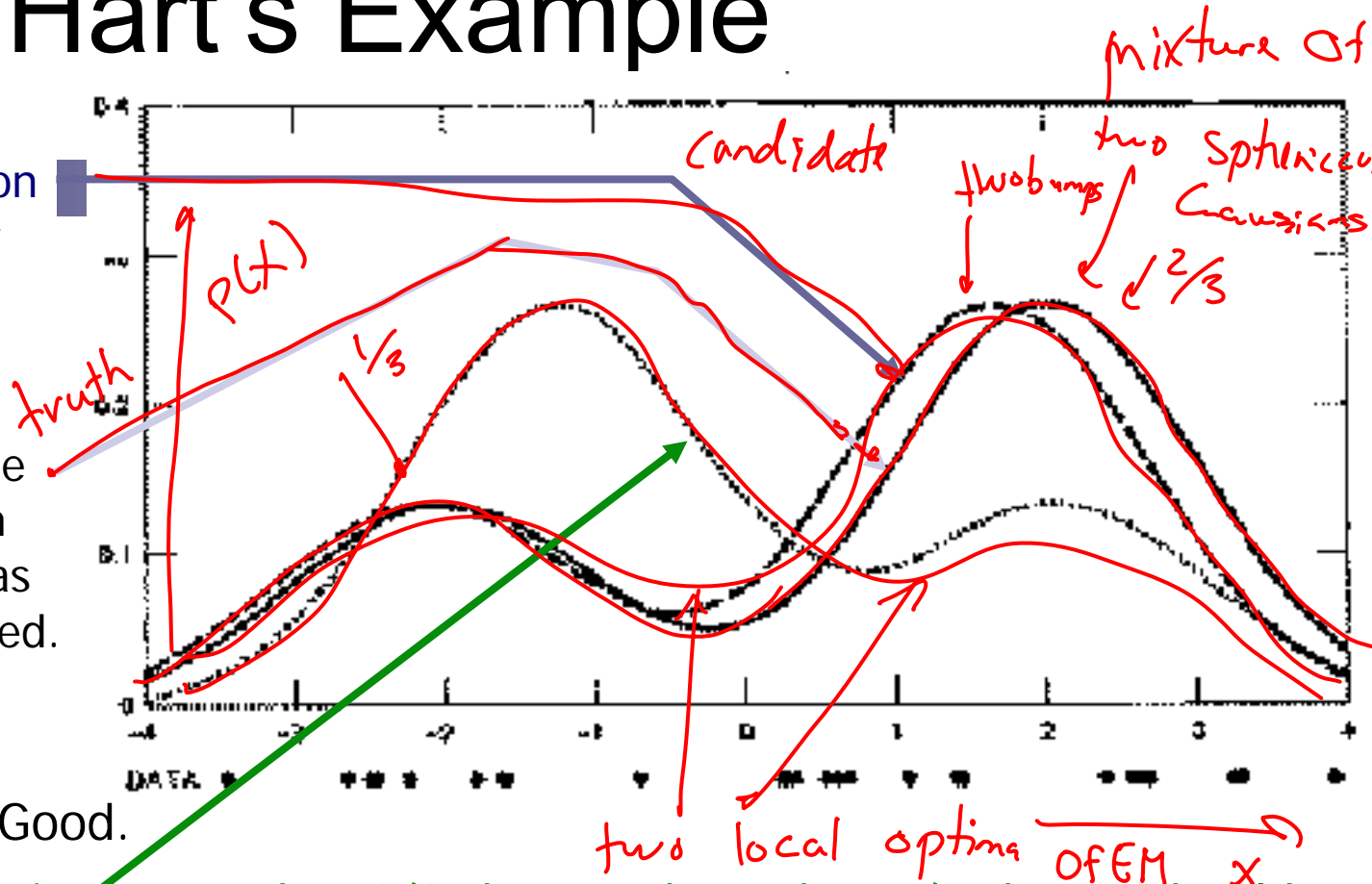
$$\begin{aligned} x_1 &= 0.608 \\ x_2 &= -1.590 \\ x_3 &= 0.235 \\ x_4 &= 3.949 \\ &\vdots \\ x_{25} &= -0.712 \end{aligned}$$



Duda & Hart's Example

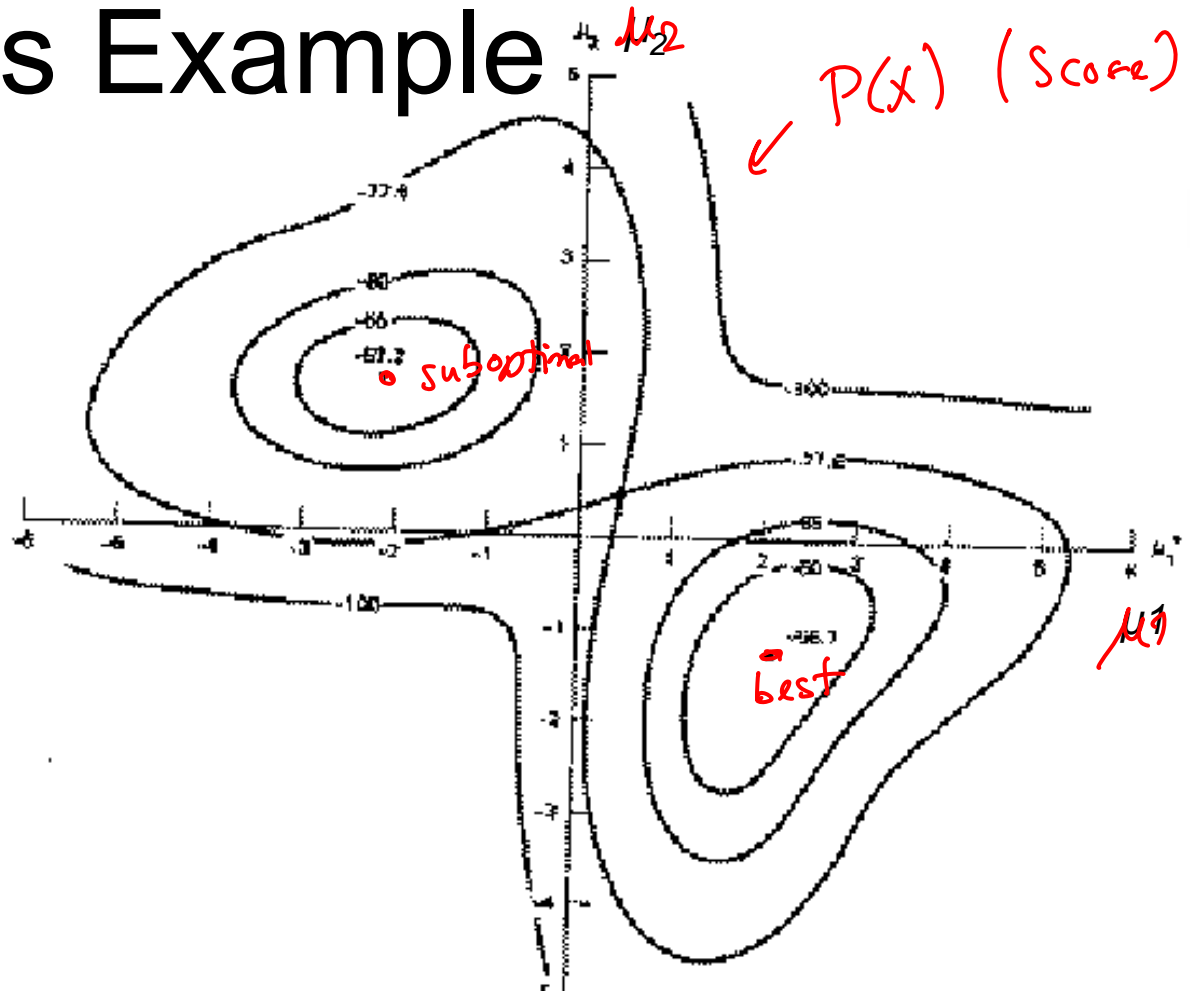
We can graph the prob. dist. function of data given our μ_1 and μ_2 estimates.

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
 $\log P(x_1, x_2 \dots x_{25} \mid \mu_1, \mu_2)$
against μ_1 (\rightarrow) and μ_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 $\mu_1, \mu_2 \dots \mu_k$

We can compute $P(\text{data} \mid \mu_1, \mu_2 \dots \mu_k)$

How do we find the μ_i 's which give max. likelihood?

- The normal max likelihood trick:

Set $\frac{\partial}{\partial \mu_j} \log \text{Prob} (\dots) = 0$

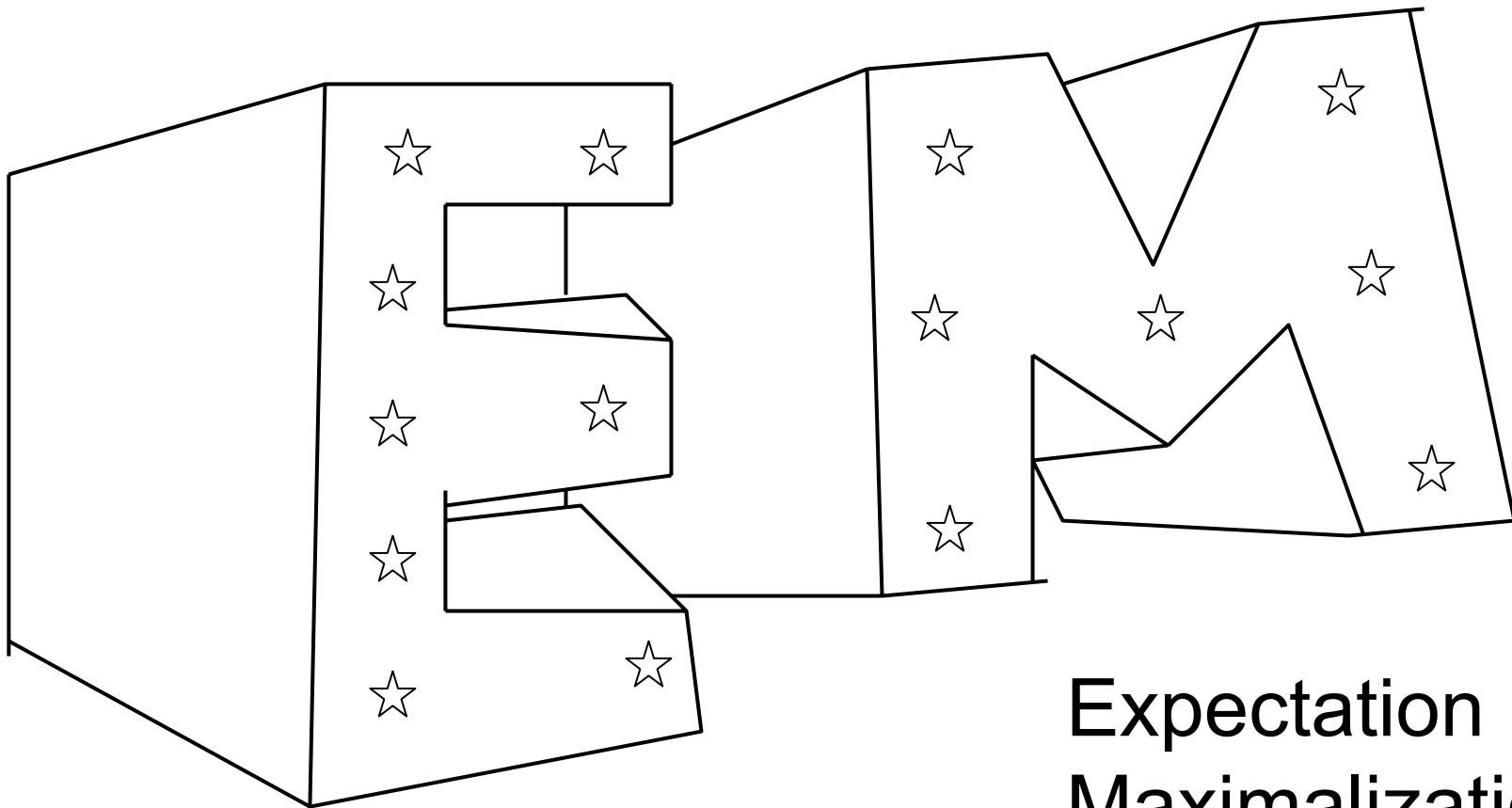
and solve for μ_i 's.

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

standard MLE: $P(A) = w_1$
 $w_1 + w_2 + w_3 + w_4$

Let events be "grades in a class"

w_1 = Gets an A

w_2 = Gets a B

w_3 = Gets a C

w_4 = Gets a D

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

$$P(B) = \mu$$

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

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

what's μ for this data?

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

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 ?

$$\hat{\mu} = \underset{\mu}{\operatorname{argmax}} p(a, b, c, d | \mu)$$
$$= \underset{\mu}{\operatorname{argmax}} \left(\frac{1}{2}\right)^a \cdot (\mu)^b \cdot (2\mu)^c \cdot \left(\frac{1}{2} - 3\mu\right)^d$$

Trivial Statistics

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

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

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

FOR MAX LIKE μ , SET $\frac{\partial \text{Log} P}{\partial \mu} = 0$

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

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

So if class got

A	B	C	D
14	6	9	10

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

~~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

*don't
tell you
how many
A's*

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 μ now?

$$\text{arg max}_{\mu} \sum_{a,b: a+b=h} p(a,b,c,d | \mu)$$

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:

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

$P(A) = \frac{1}{2}$
 $P(B) = \frac{1}{4}$

$h = 10$ $\mu = \frac{1}{4}$

$a = \frac{\frac{1}{2}}{\frac{1}{2} + \frac{1}{4}} \cdot 10 = \frac{2}{3} \cdot 10$

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

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

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

MAXIMIZATION

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

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

E.M. for our Trivial Problem

REMEMBER

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

$$P(B) = \mu$$

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

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

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

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

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

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

= max like est. of μ given $b^{(t)}$

expected number of b's

E-step

M-step

best μ for this # of b's.

Continue iterating until converged.

Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.

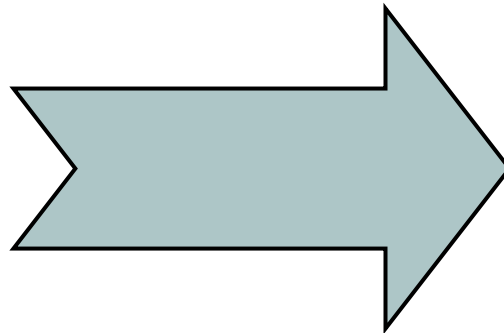
E.M. Convergence

- 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]

$p(\text{data} \mid \mu)$ Converge!

In our example,
suppose we had

$$\begin{aligned} h &= 20 \\ c &= 10 \\ d &= 10 \\ \mu^{(0)} &= 0 \end{aligned}$$



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

3.187
out of
20

Back to Unsupervised Learning of GMMs – a simple case

Remember:

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

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

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

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

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

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

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 μ_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 datapoints 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*

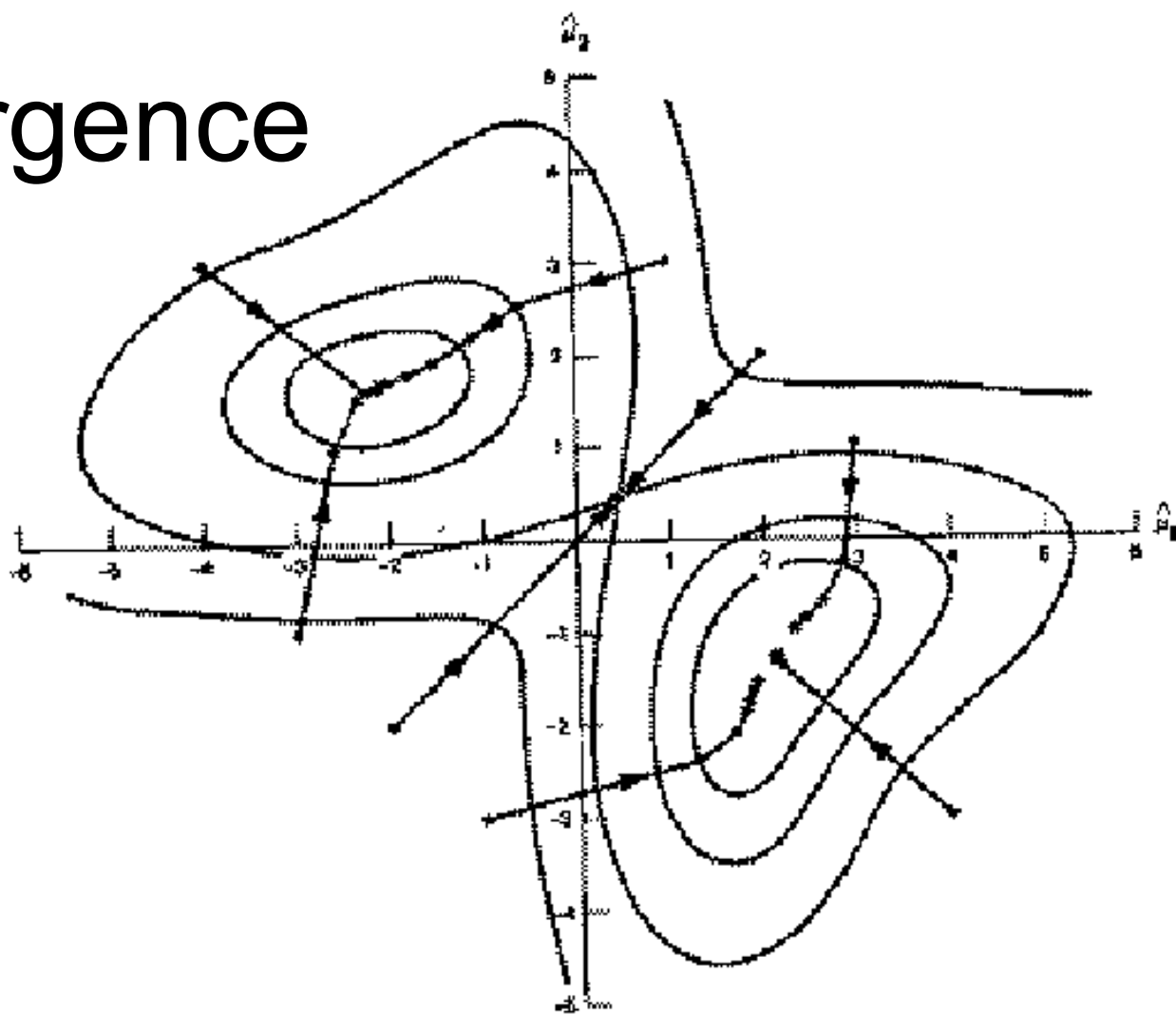
M-step

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

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

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

$p_i^{(t)}$ 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

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

Just evaluate a Gaussian at x_j

M-step

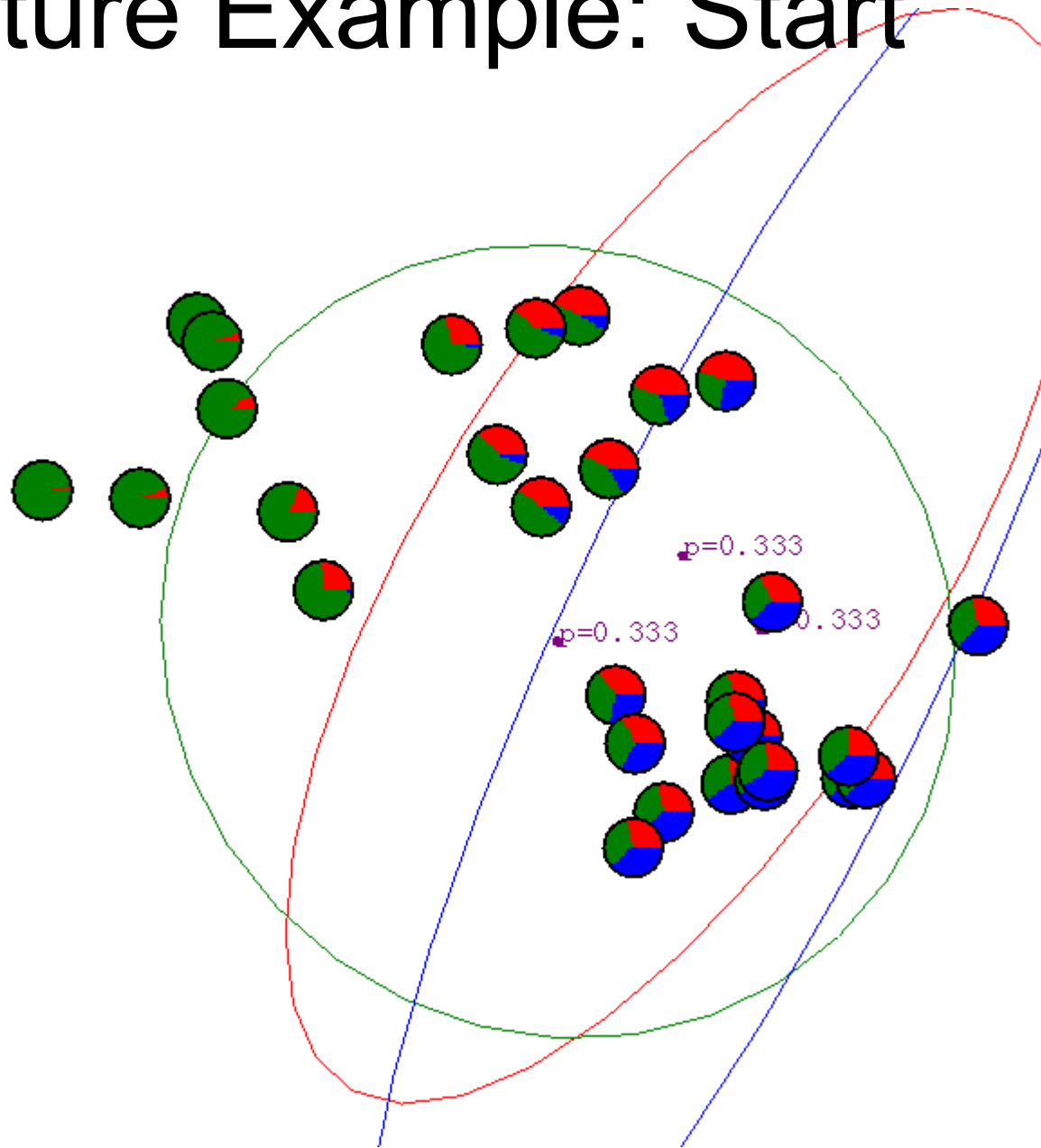
Compute Max. like μ given our data's class membership distributions

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

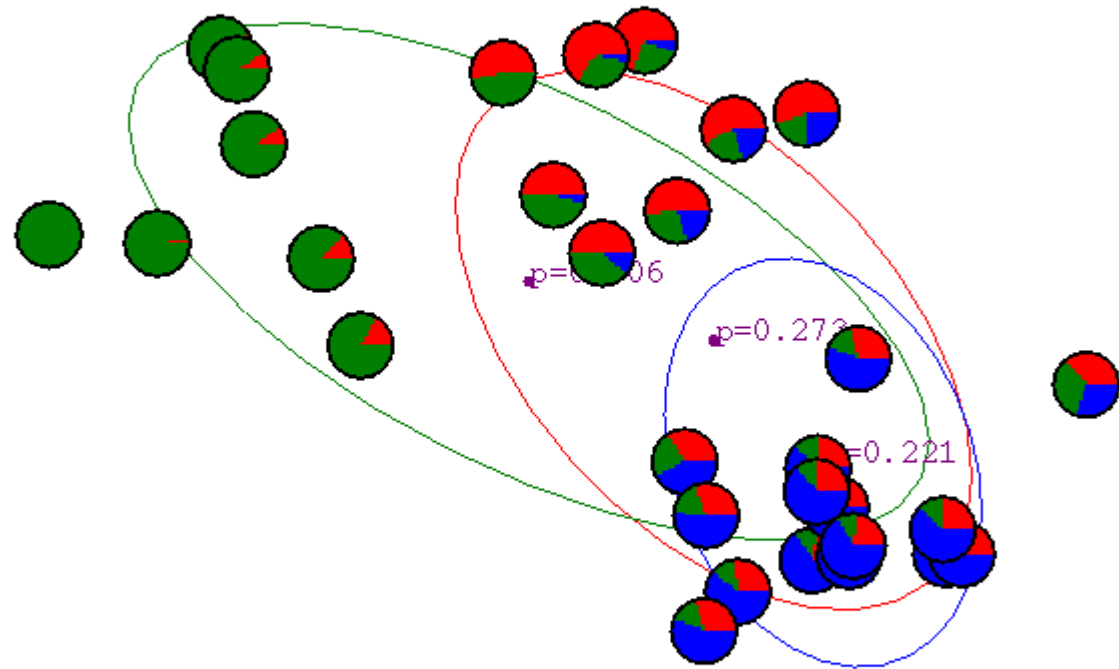
$$p_i^{(t+1)} = \frac{\sum_j P(y = i | x_j, \lambda_t)}{m}$$

$m = \#records$

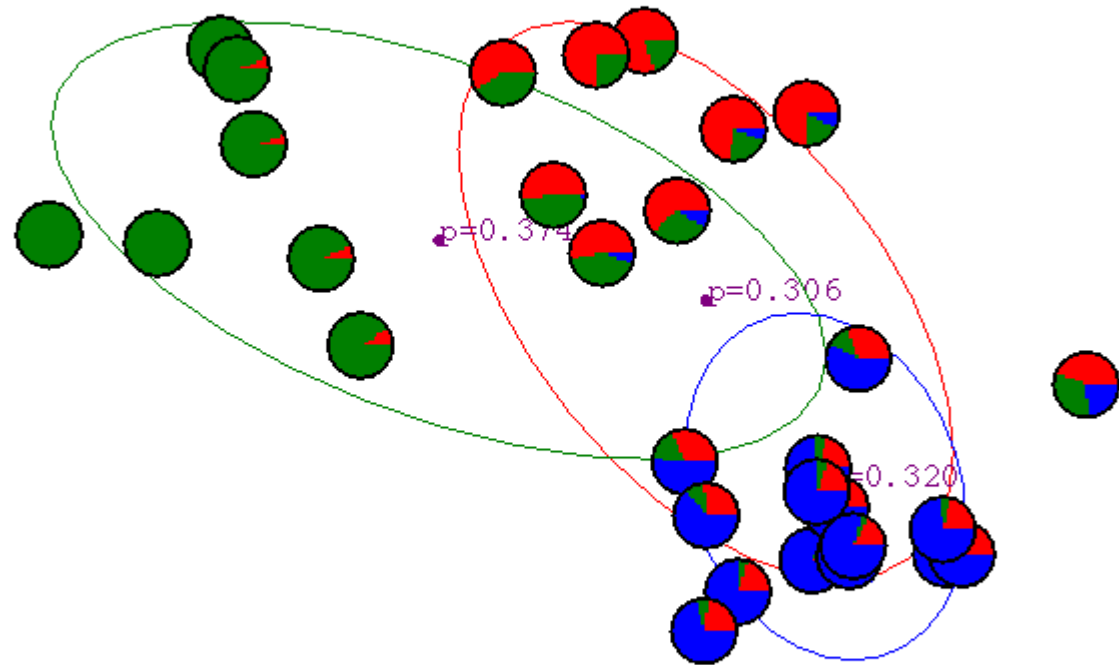
Gaussian Mixture Example: Start



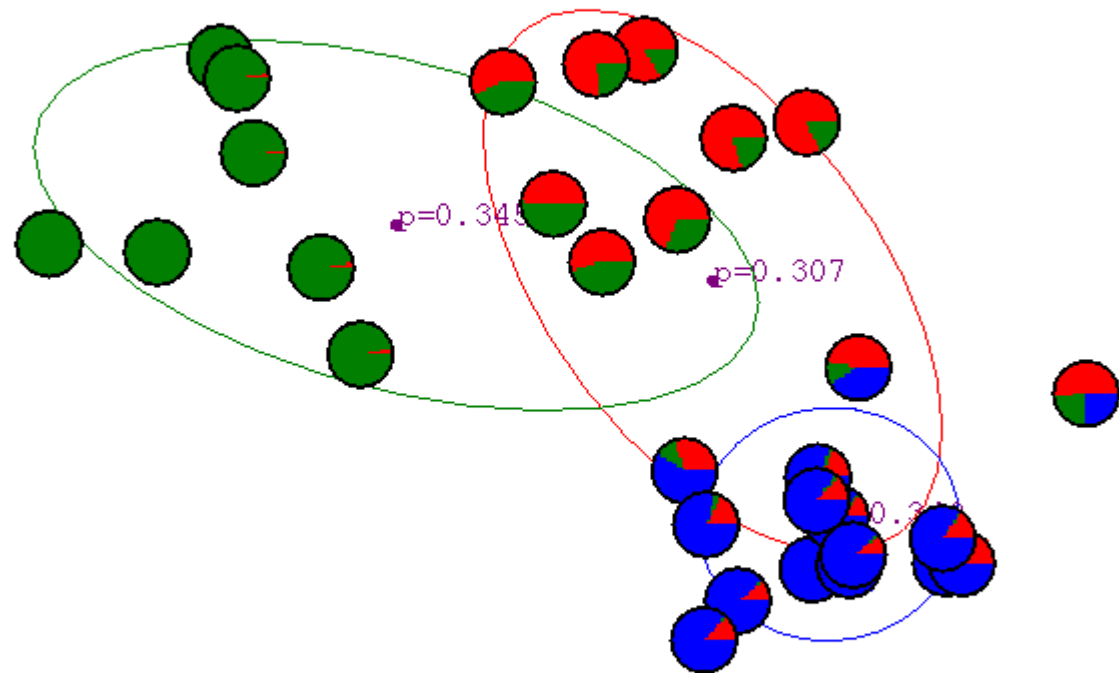
After first iteration



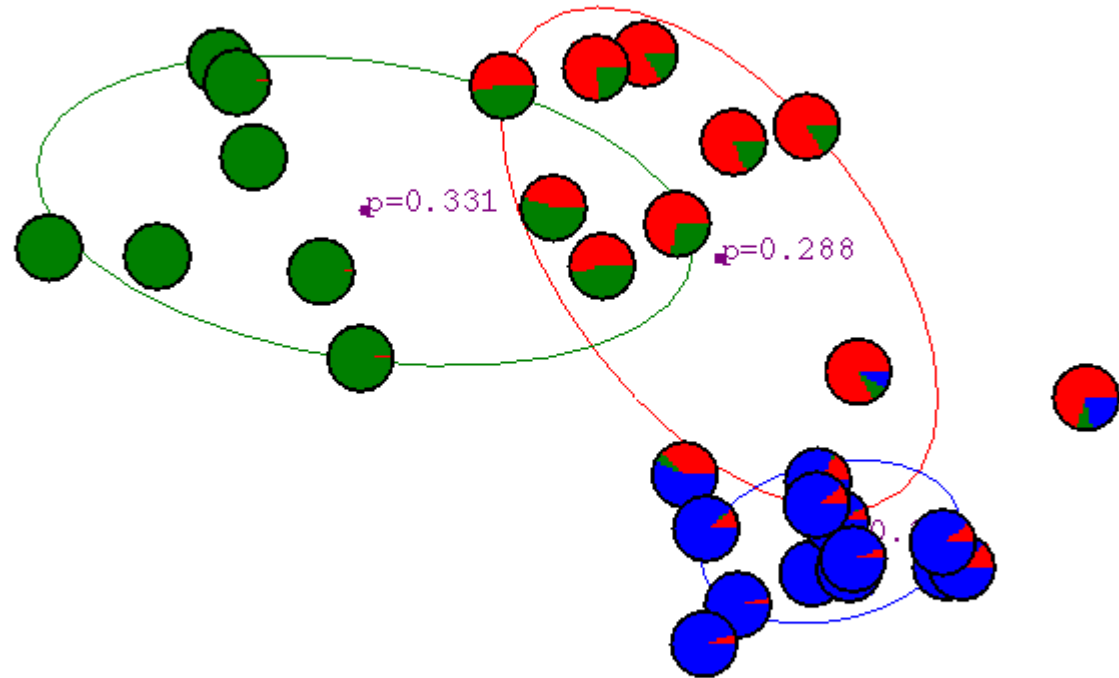
After 2nd iteration



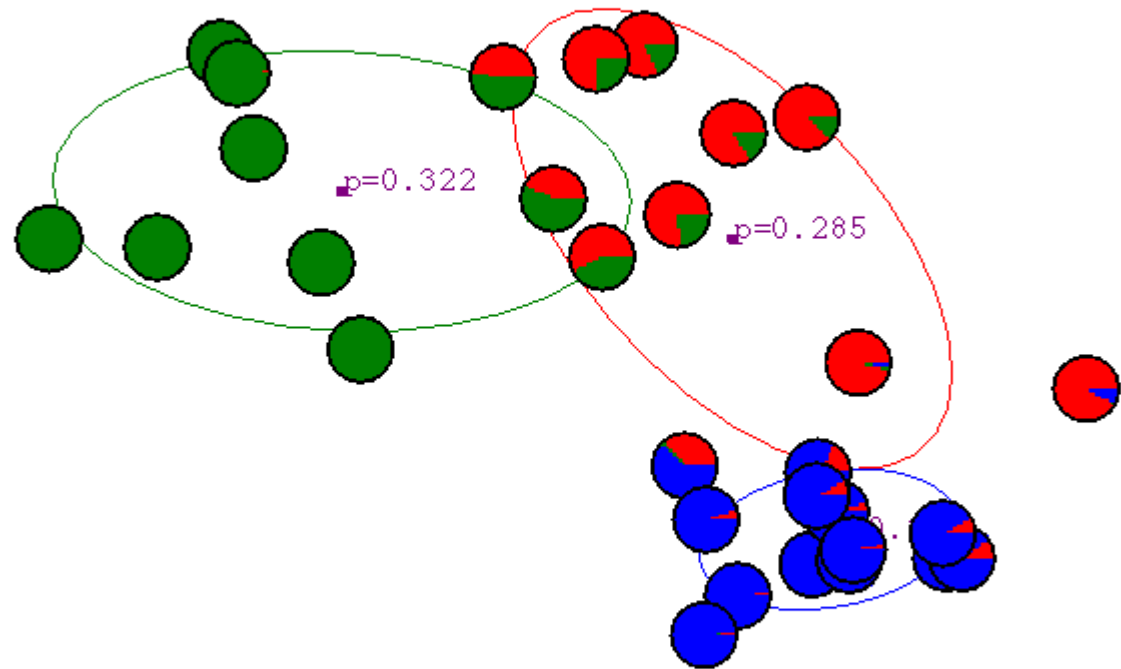
After 3rd iteration



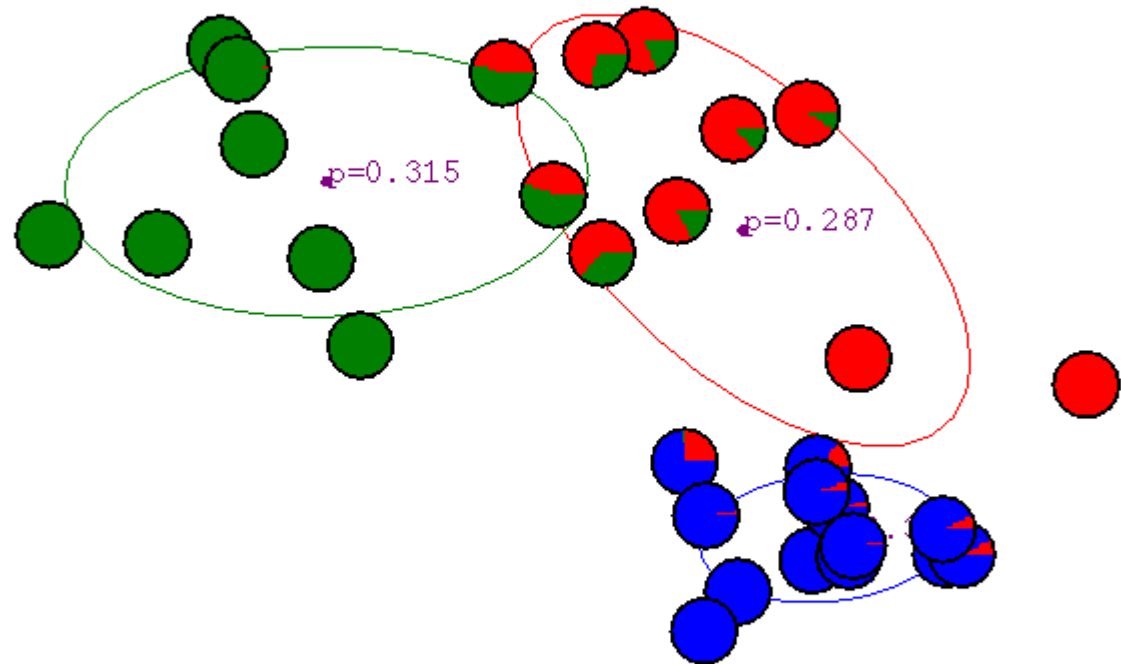
After 4th iteration



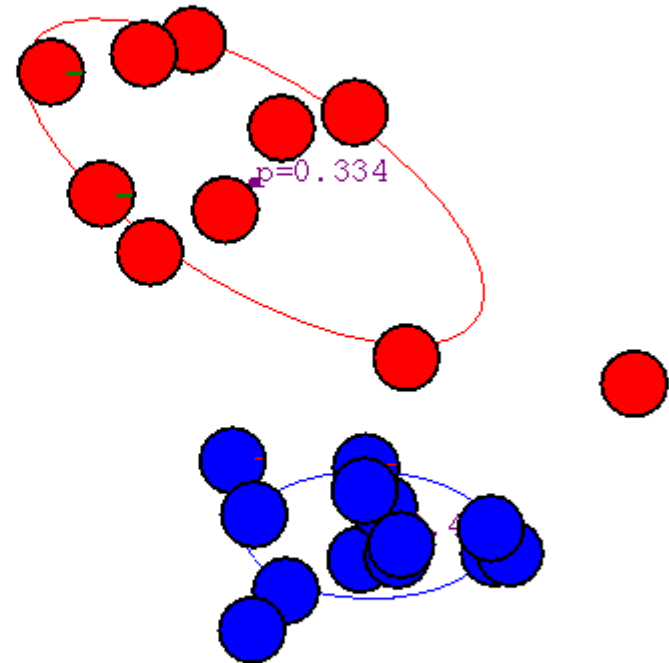
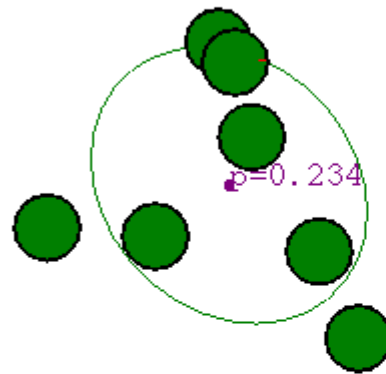
After 5th iteration



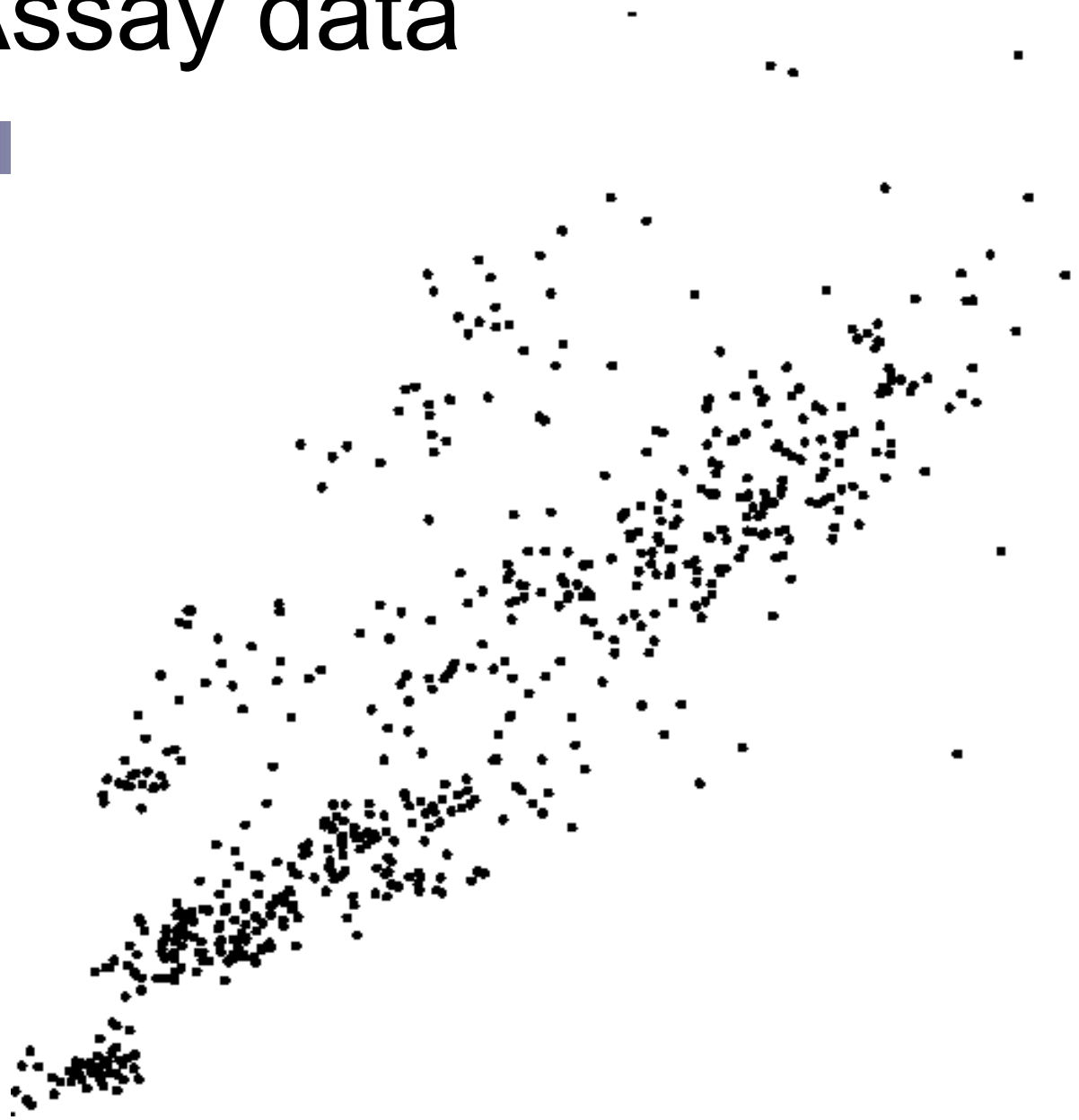
After 6th iteration



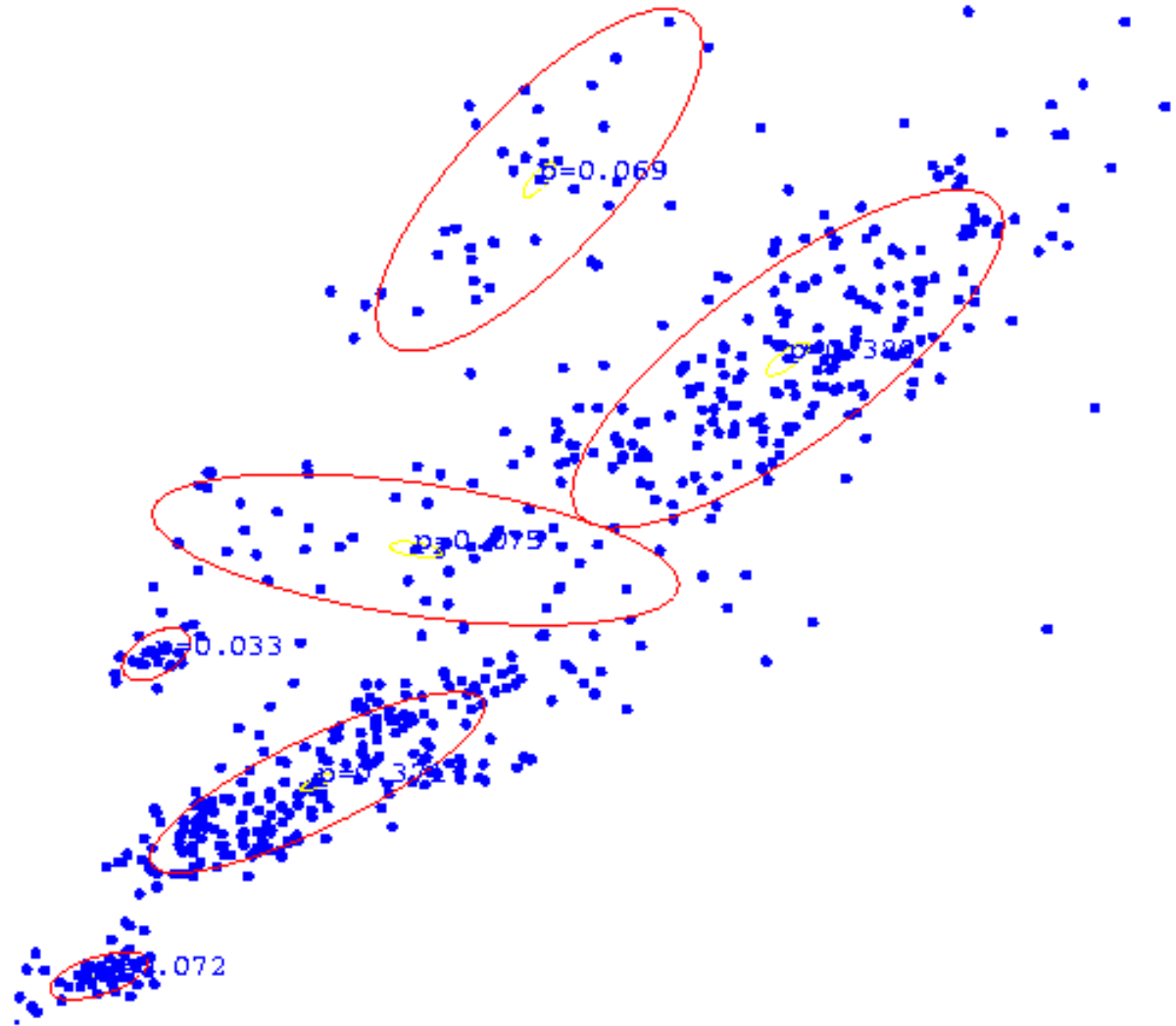
After 20th iteration




Some Bio Assay data



GMM clustering of the assay data





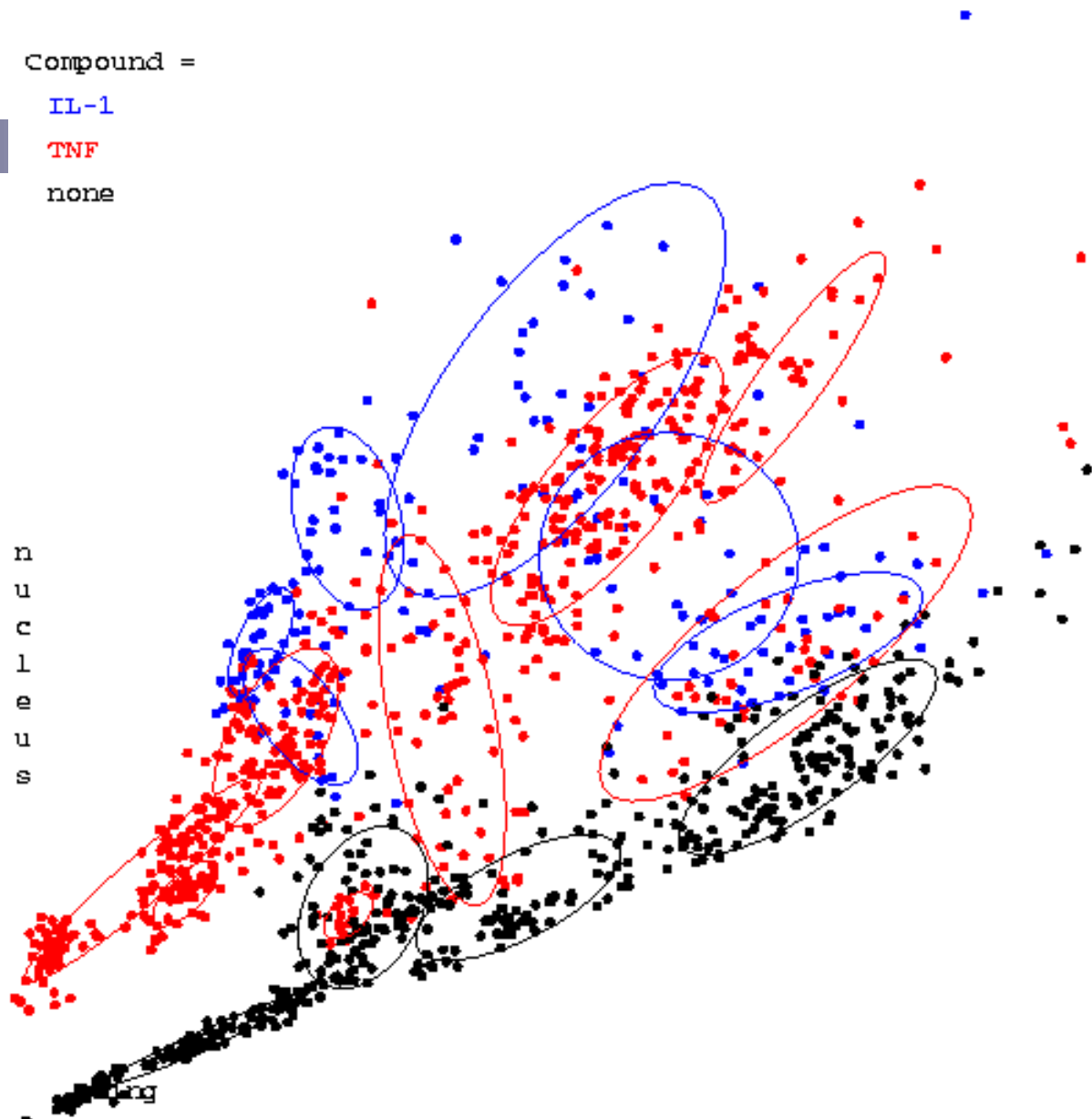
Resulting Density Estimator



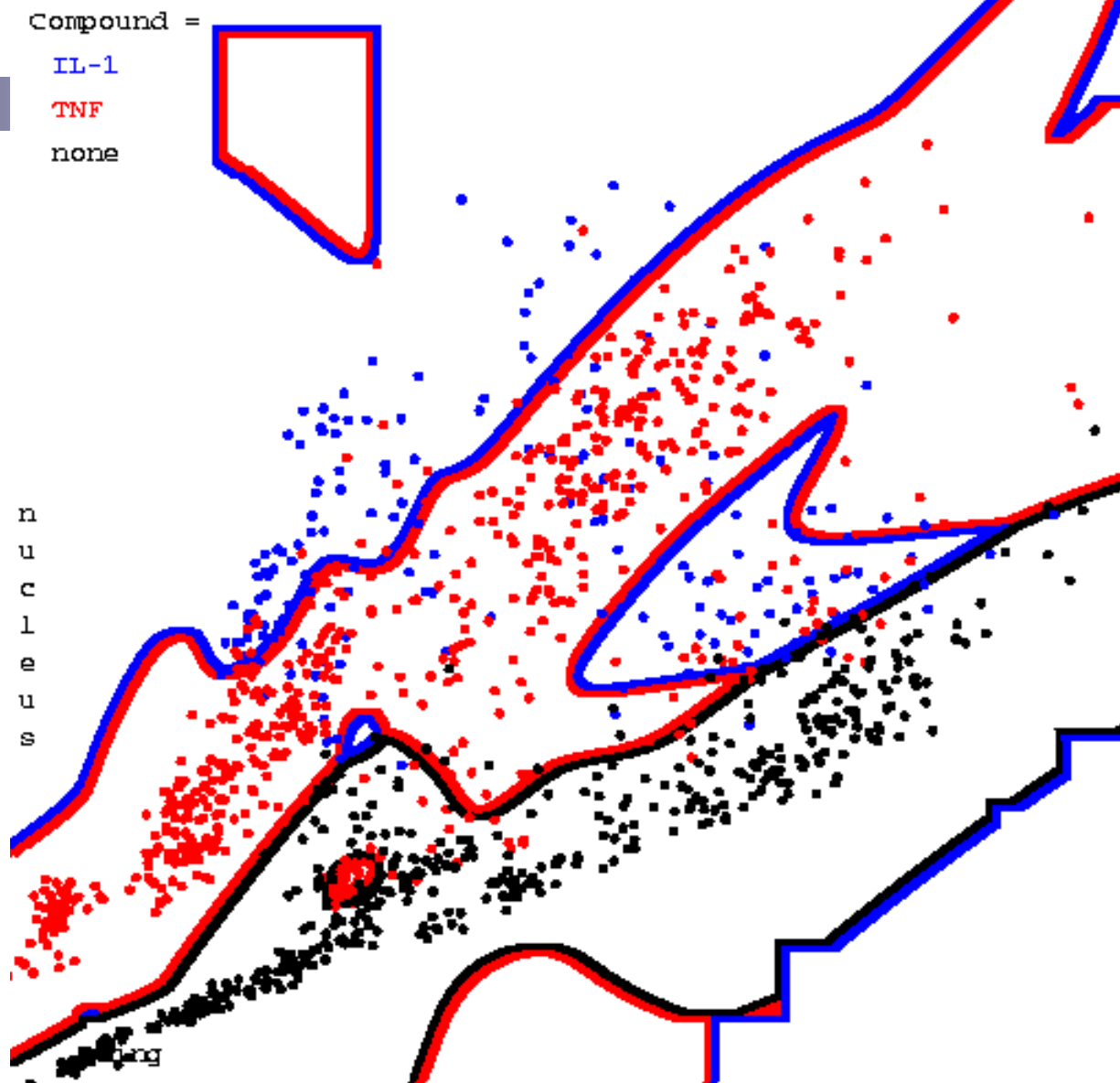
Three classes of assay

(each learned with its own mixture model)

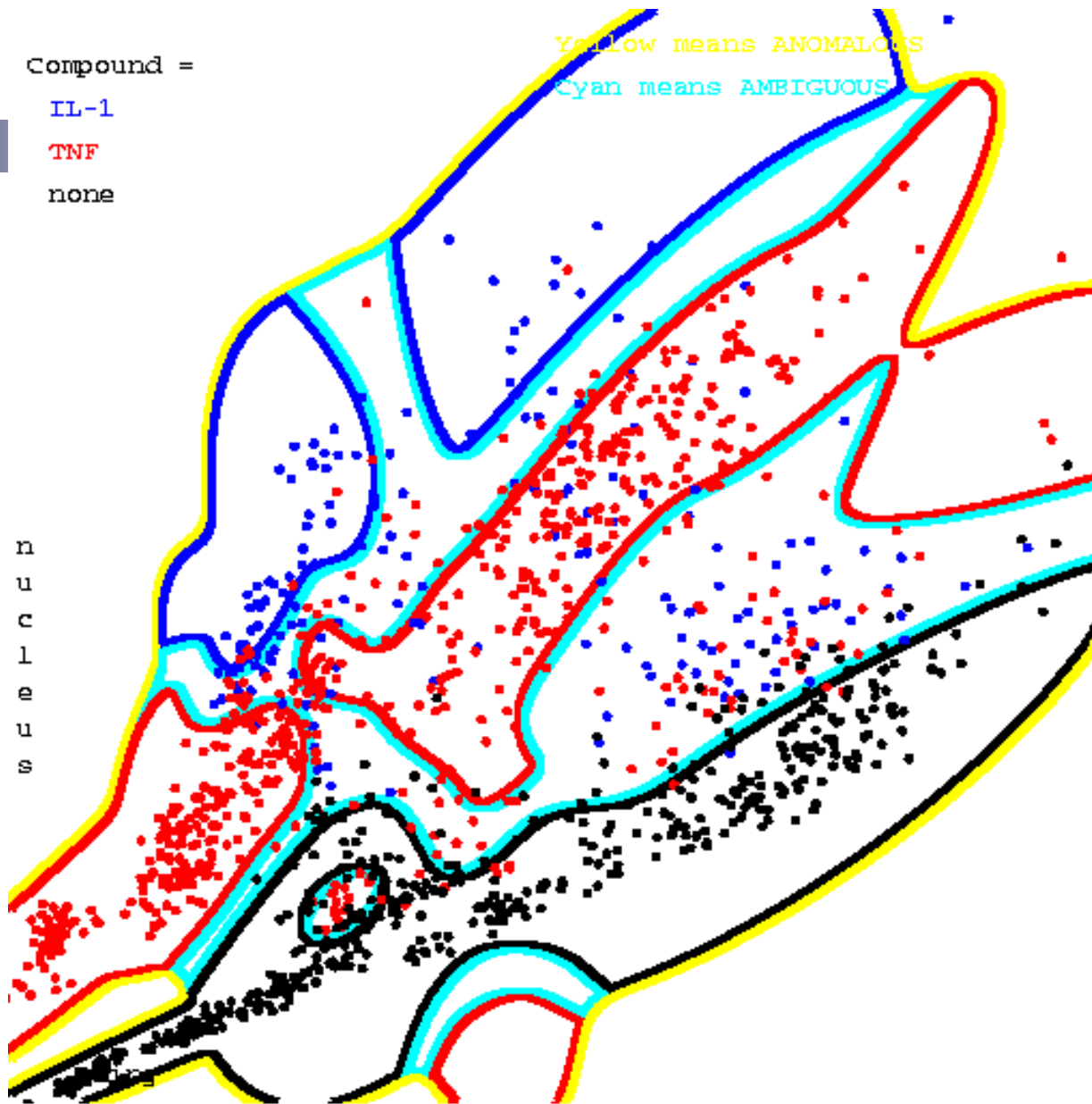
Compound =
IL-1
TNF
none



Resulting Bayes Classifier



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



Yellow means
anomalous

Cyan means
ambiguous

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 DOES

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>