

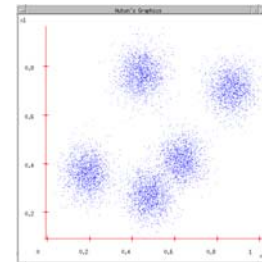
# Data Structures for Fast Mixture Models

The Auton Lab  
Carnegie Mellon University



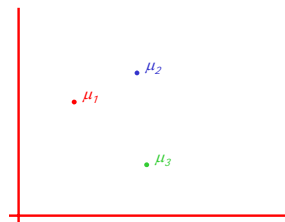
[www.autonlab.org](http://www.autonlab.org)

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



## The GMM assumption

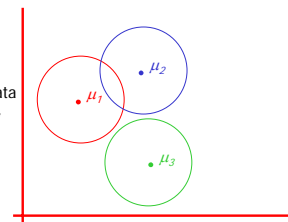
- There are  $k$  components. The  $i$ 'th component is called  $\omega_i$
- Component  $\omega_i$  has an associated mean vector  $\mu_i$



## 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:

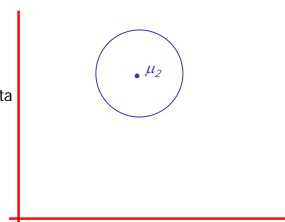


## The GMM assumption

- 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:

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

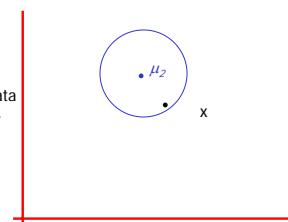


## The GMM assumption

- 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:

- Pick a component at random. Choose component  $i$  with probability  $P(\omega_i)$ .
- Datapoint  $\sim N(\mu_i, \sigma^2 I)$

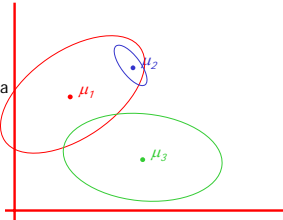


## The General GMM assumption

- 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  $\Sigma_i$

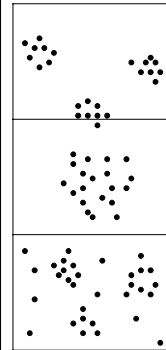
Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component  $i$  with probability  $P(\omega_i)$ .
2. Datapoint  $\sim N(\mu_i, \Sigma_i)$



7

## Unsupervised Learning: not as hard as it looks



Sometimes easy

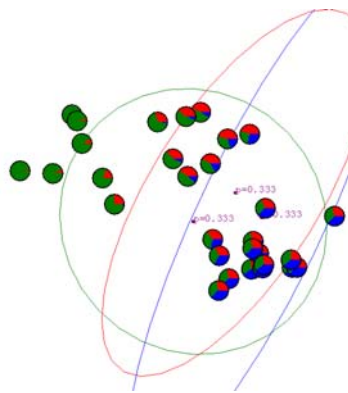
Sometimes impossible

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

8

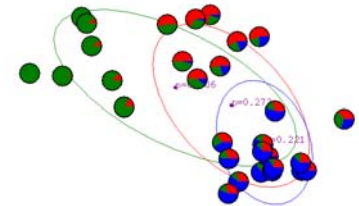
## Gaussian Mixture Example: Start



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

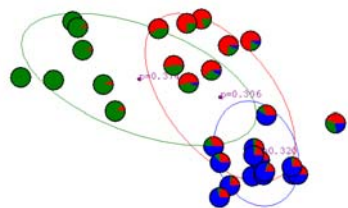
9

## After first iteration



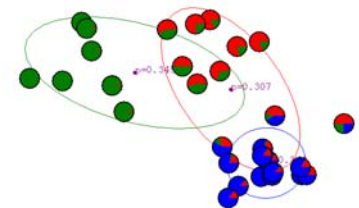
10

## After 2nd iteration



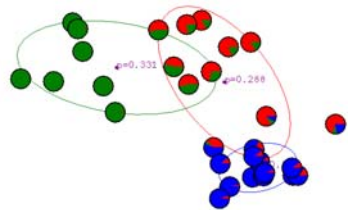
11

## After 3rd iteration



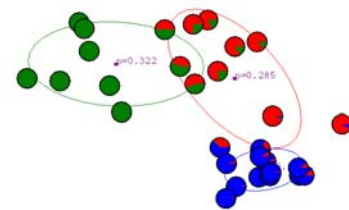
12

After 4th  
iteration



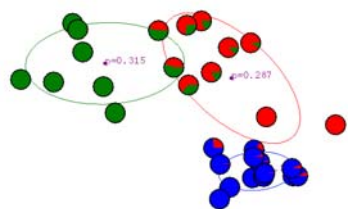
13

After 5th  
iteration



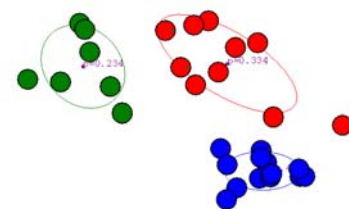
14

After 6th  
iteration



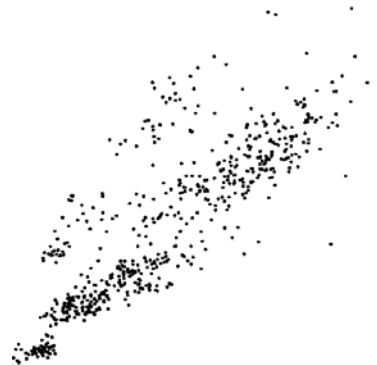
15

After 20th  
iteration



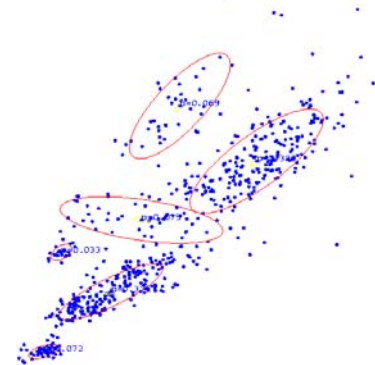
16

Some Bio  
Assay  
data



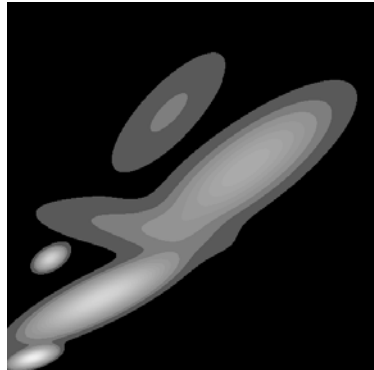
17

GMM  
clustering  
of the  
assay data



18

## Resulting Density Estimator



19

## The MRKD-tree approach

Kdtrees are a multidimensional data structure for representing sets of points in k-dimensional space. Invented in the early 70's by Jon Bentley.

MRKD-trees are Kd-trees containing extra information (cached sufficient statistics) including covariance matrices at every node in order to speed up statistical queries. [Deng and Moore, 1995].

With MRKD-trees we have been able to speed up a wide range of multidimensional statistics and machine learning algorithms [Deng and Moore 95, Moore Schneider and Deng 97, Moore 1998, Pelleg and Moore 1999]

Turns out to be very useful for clustering. Typically, instead of visiting millions of datapoints at each iteration, simply visit a few thousand tree nodes. Some quadratic programming and algebraic bounds on Bayes Rule guarantees correct performance.

20

## Computing likelihood of datapoints...



suppose you want to compute the sum of log-likelihoods of all the blue dots given they'd been generated by the big red Gaussian.

## Computing likelihood of datapoints...

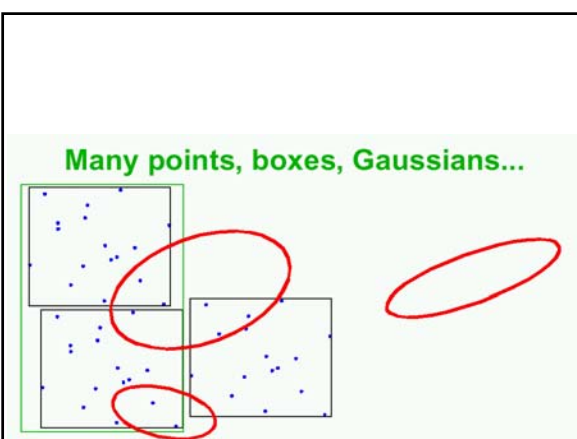


suppose we happen to know their bounding box

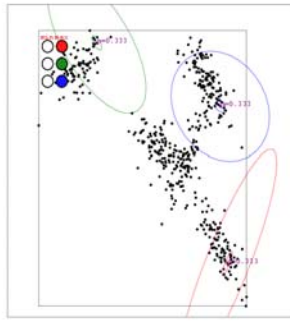
## Computing likelihood of datapoints...



Without visiting the points individually, we can put bounds on their contributions to the Gaussian. Sometimes those bounds'll be tight enough...

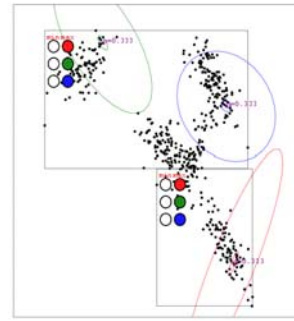


**MRKD at the top node**



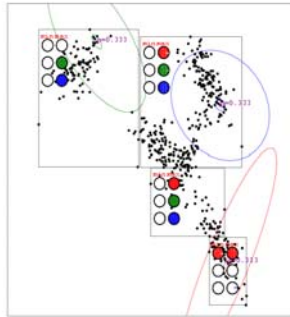
25

**MRKD at level one nodes**



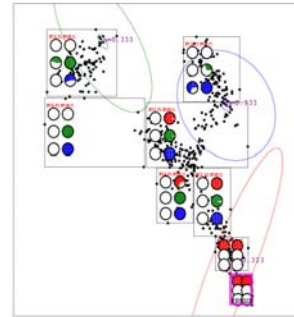
26

**MRKD at level two nodes**



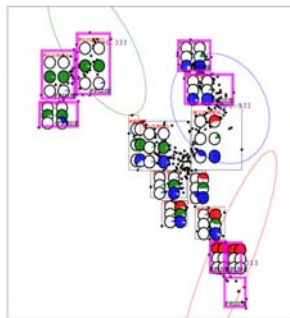
27

**MRKD at level three nodes**



28

**MRKD at level four nodes**



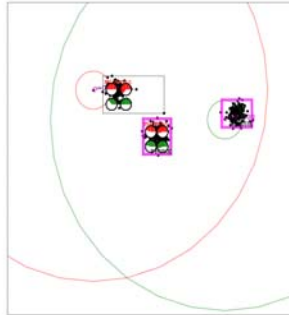
29

**MRKD at level five nodes**



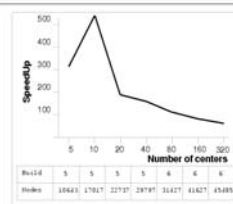
30

### Another MRKD prune example



31

### Effect of Number of Centers

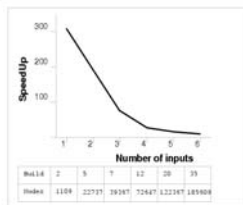


When  $k=320$ : SlowSecs = 9278, FastSecs = 143.3, LLDiff = 0.004

Conventional EM slows down linearly with the number of classes. Fast EM is clearly sublinear, with a 70-fold speedup even with 320 classes. Note how the tree size grows. This is because more classes mean a more uniform data distribution and fewer datapoints "sharing" tree leaves.

32

### Effect of Dimensionality

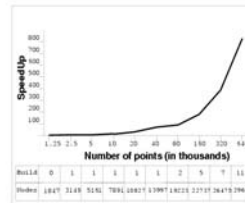


When  $M=6$ : SlowSecs = 2742, FastSecs = 310.25, LLDiff = 0.08

As with many kd-tree algorithms, the benefits decline as dimensionality increases, yet even in 6 dimensions, there is an 8-fold advantage.

33

### Effect of Number of Datapoints

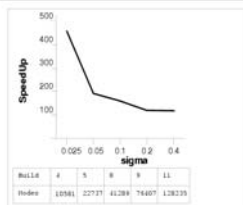


When  $R=640,000$ : SlowSecs = 2385, FastSecs = 3, LLDiff = 0.009

As  $R$  increases so does the computational advantage, essentially linearly. The tree-build time (11 seconds at worst) is a tiny cost compared with even just one iteration of Regular EM (2385 seconds, on the big dataset.)

34

### Effect of Initial Gaussian Widths

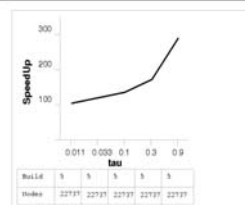


When  $\text{SigmaStart} = 0.4$ : SlowSecs = 583, FastSecs = 4.75, LLDiff = 0.003

Even with very wide Gaussians, with wide support, we still get large savings. The nodes that are pruned in these cases are rarely nodes with one class owning all the probability, but instead are nodes where all classes have non-zero, but little varying, probability.

35

### Effect of TAU (the willingness to prune)



When  $\text{Tau}=0.9$ : SlowSecs = 584.5, FastSecs = 2, LLDiff = 0.06

The larger tau, the more willing we are to prune during the tree search, and thus the faster we search, but the less accurately we mirror EM's statistical behavior.

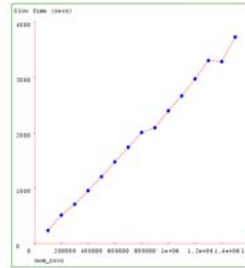
36

### Nodes visited during an EM pass



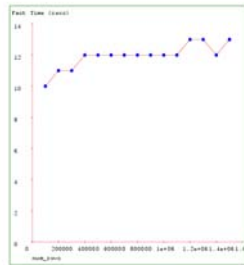
27

### Sky Survey Data: Time Taken by the Slow method.



28

### Sky Survey Data: Time Taken by the Fast method



29