

Mixtures of Gaussians continued

Machine Learning – CSE446

Carlos Guestrin

University of Washington

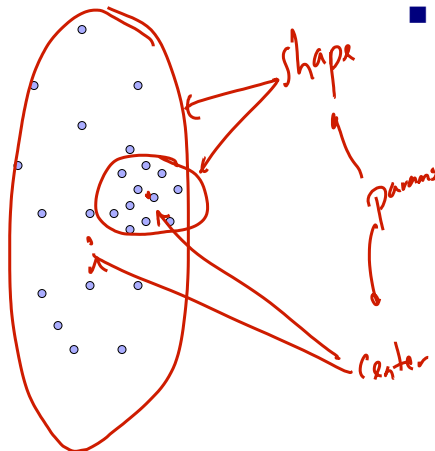
May 17, 2013

©Carlos Guestrin 2005-2013

1

(One) bad case for k-means

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



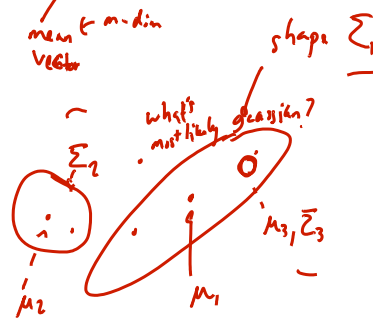
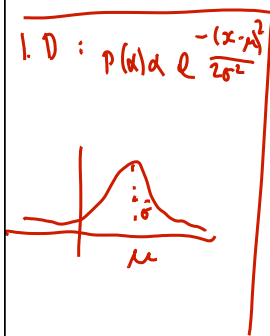
©Carlos Guestrin 2005-2013

2

Gaussians in m Dimensions

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

m-dim (pointing to Σ)
Covariance matrix (pointing to Σ)
 Σ is $n \times n$ (pointing to Σ)



©Carlos Guestrin 2005-2013

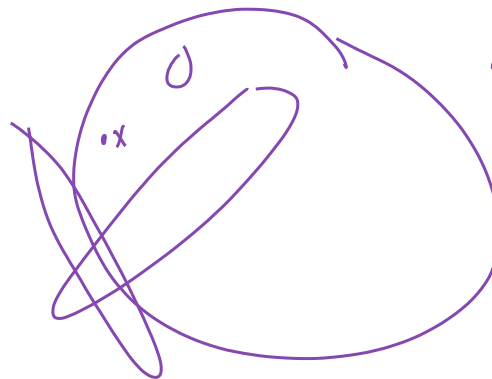
3

Suppose You Have a Gaussian For Each Class

$$\|a-b\|_{\Sigma} = (a-b)^T \Sigma^{-1} (a-b)$$

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

$$p(y=i|x)$$



a distance metric per class

©Carlos Guestrin 2005-2013

4

Gaussian Bayes Classifier

- You have a Gaussian over \mathbf{x} for each class $y=i$:

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

- But you need probability of class $y=i$ given \mathbf{x} :

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

- Thank you Bayes Rule!!

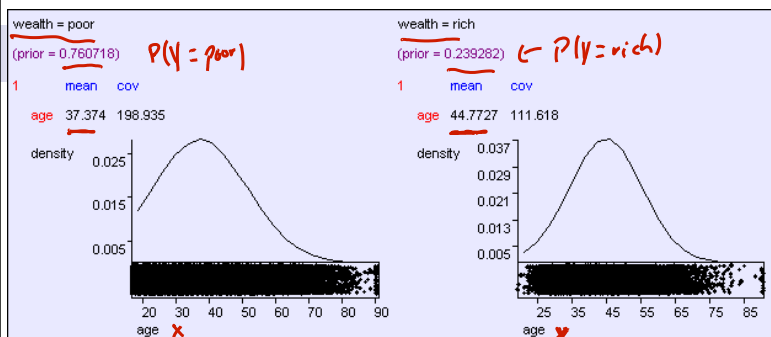
$$P(y = i | \mathbf{x}) = \frac{P(\mathbf{x} | y = i)P(y = i)}{p(\mathbf{x})}$$

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

©Carlos Guestrin 2005-2013

5

Predicting wealth from age



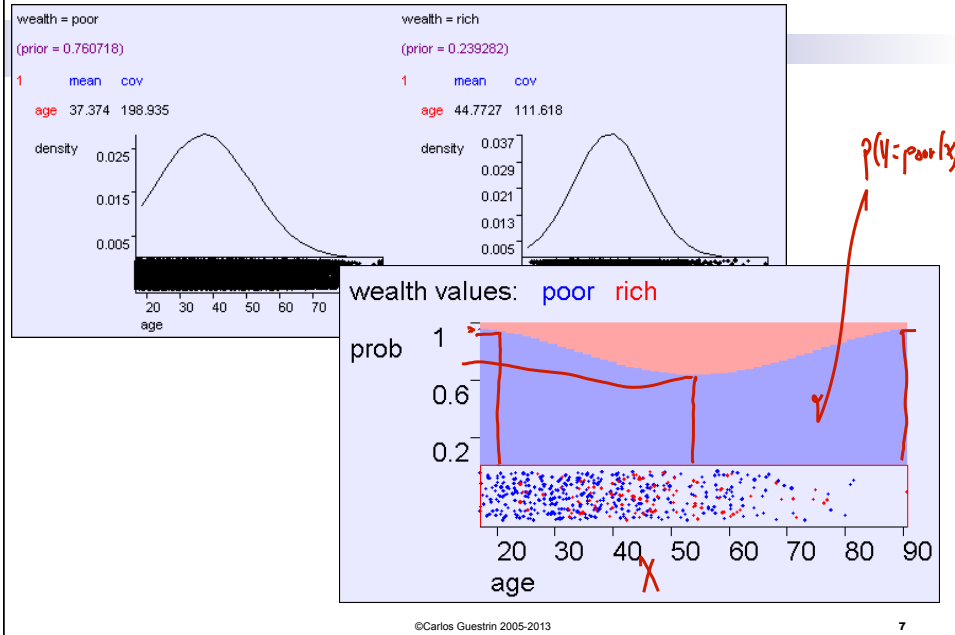
$$P(y = \text{poor} | \mathbf{x})$$

$$\propto P(y = \text{poor}) \cdot p(\mathbf{x} | y = \text{poor})$$

©Carlos Guestrin 2005-2013

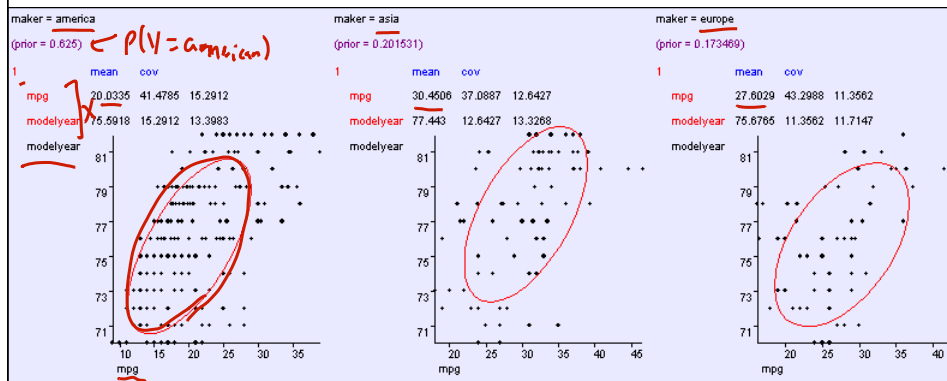
6

Predicting wealth from age



Learning modelyear ,
mpg ---> maker

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma_2^2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma_m^2 \end{pmatrix}$$

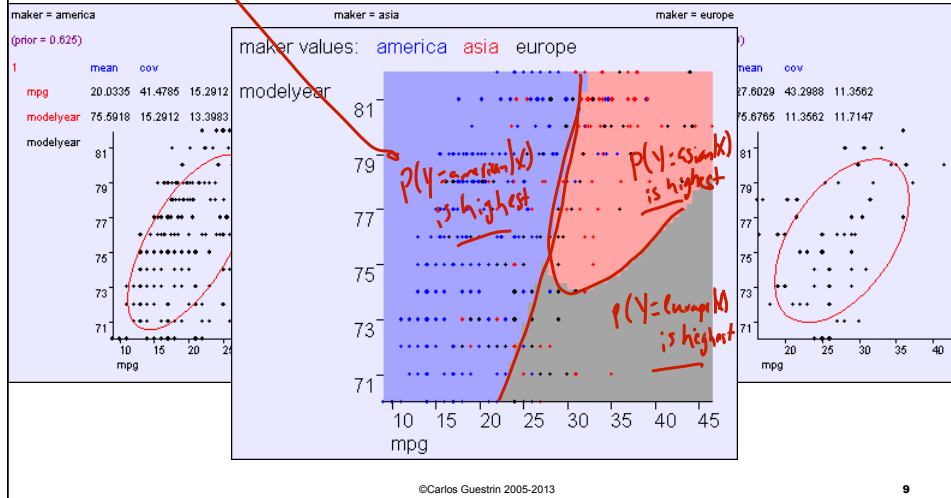


General: $O(m^2)$

parameters

\uparrow m dims of x
 (classification $\arg\max_y P(Y=y|x)$)

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{12} & \sigma_2^2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \cdots & \sigma_m^2 \end{pmatrix}$$



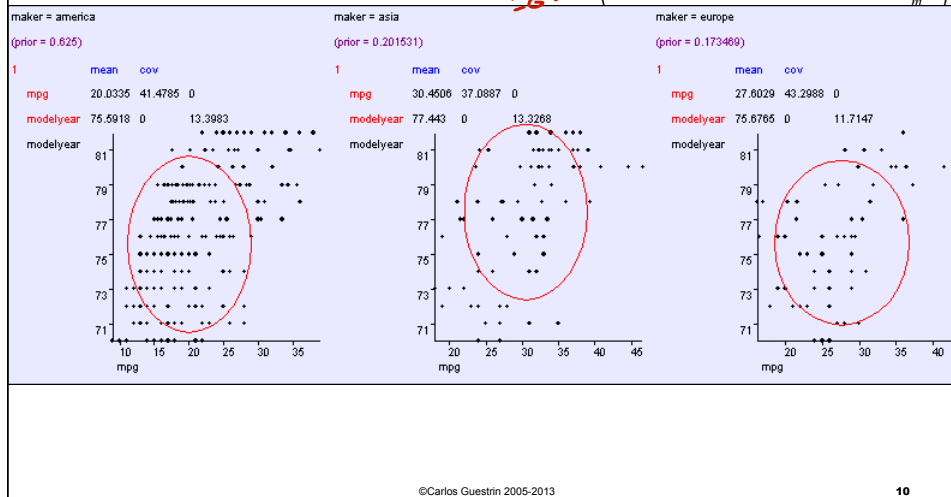
Aligned: $O(m)$

parameters

free to scale

$\Sigma =$
 off-diagonal

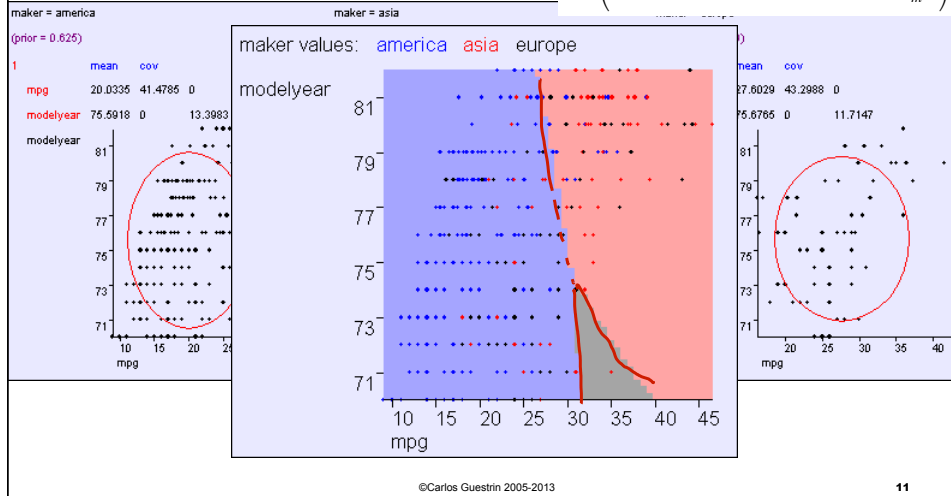
$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma_2^2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma_3^2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma_{m-1}^2 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma_m^2 \end{pmatrix}$$



Aligned: $O(m)$
parameters

$\Sigma =$

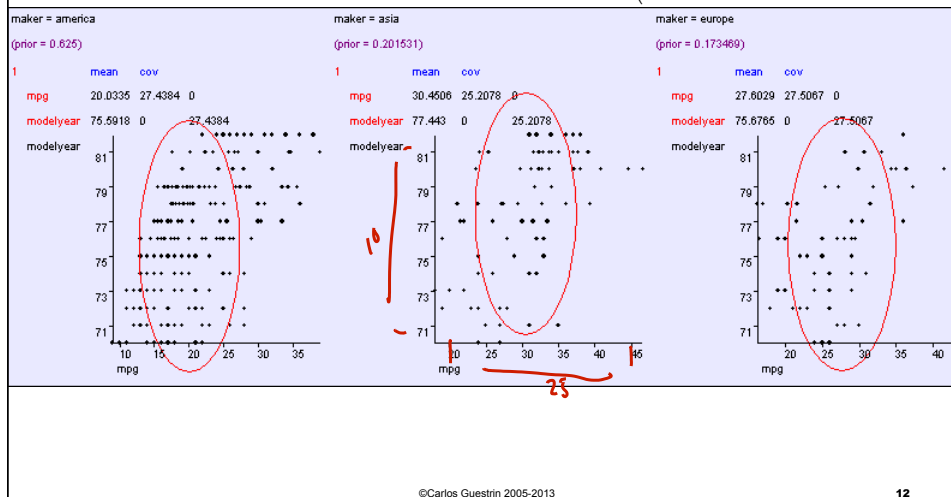
$$\begin{pmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_2^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_3^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{m-1}^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_m^2 \end{pmatrix}$$



Spherical: $O(1)$
cov parameters

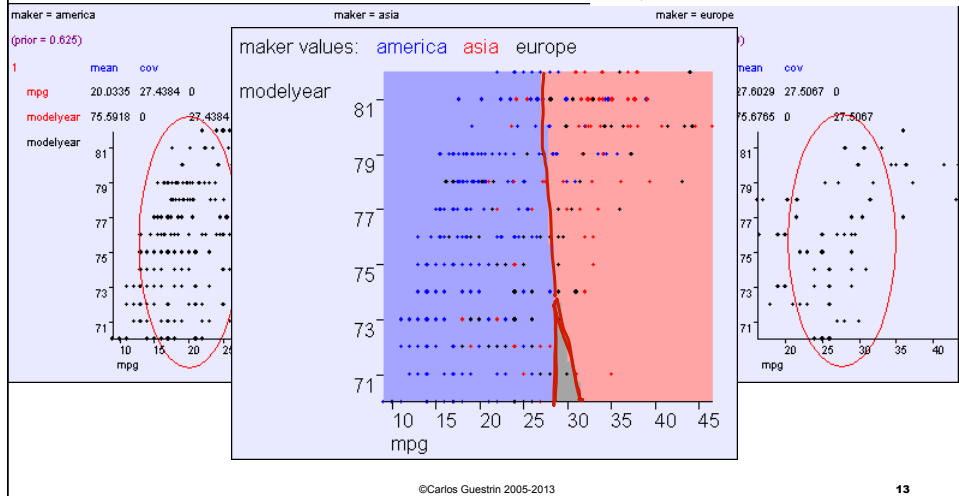
all diag
entries
= σ^2

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



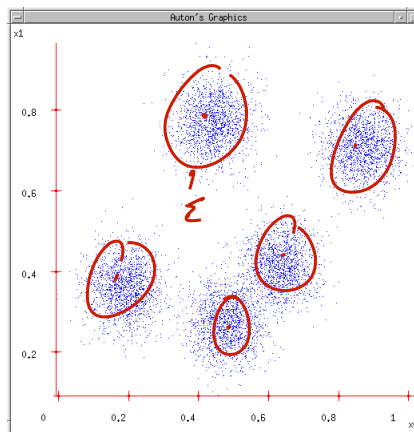
Spherical: $O(1)$
cov parameters

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



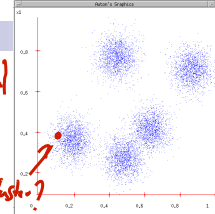
Next... back to Density Estimation

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



But we don't see class labels!!!

- MLE: in classification $\mathbf{x}^j = \{\text{MPG}, \text{Year}\}$, $y^j = \{E, Am, Asian\}$
 - $\arg\max \prod_j P(y^j, \mathbf{x}^j)$



- But we don't know y^j !!!
- Maximize marginal likelihood:
 - $\arg\max \prod_j P(\mathbf{x}^j) = \arg\max \prod_j \sum_{i=1}^k P(y^j=i, \mathbf{x}^j)$

sum or average out
unknowns/
unobserved variables

$$P(y=Am, \mathbf{x}^j) + P(y=Asian, \mathbf{x}^j) + P(y=E, \mathbf{x}^j)$$

$$P(a,b) = P(a \text{ and } b) \\ = P(Y=E, X=\{22, 1975\})$$

©Carlos Guestrin 2005-2013

15

Special case: spherical Gaussians and hard assignments

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

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

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

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

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

- Same as K-means!!!

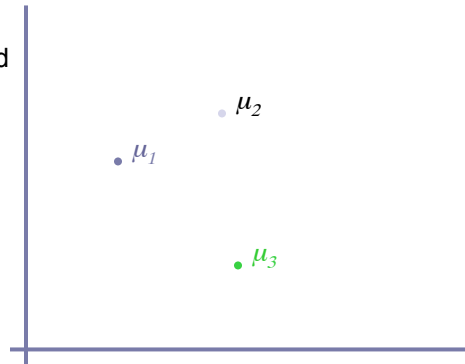
$$\max_{\mu, \sigma} F(\mu, \sigma) = \max_{\mu, \sigma} \ln F(\mu, \sigma) = \max_{\mu, \sigma} \sum_{j=1}^N -\frac{1}{2\sigma^2} \|\mathbf{x}^j - \mu_{C(j)}\|^2 \\ = \min_{\mu, \sigma} \sum_{j=1}^N \|\mathbf{x}^j - \mu_{C(j)}\|^2$$

©Carlos Guestrin 2005-2013

16

The GMM assumption

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



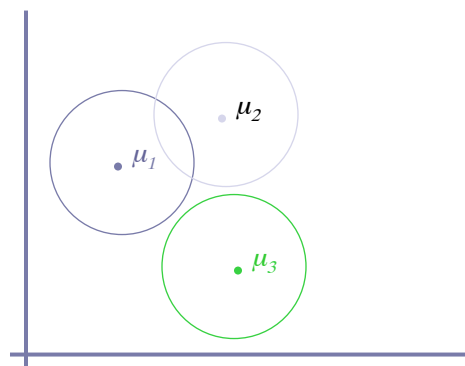
©Carlos Guestrin 2005-2013

17

The GMM assumption, in spherical case

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

Each data point is generated according to the following recipe:



©Carlos Guestrin 2005-2013

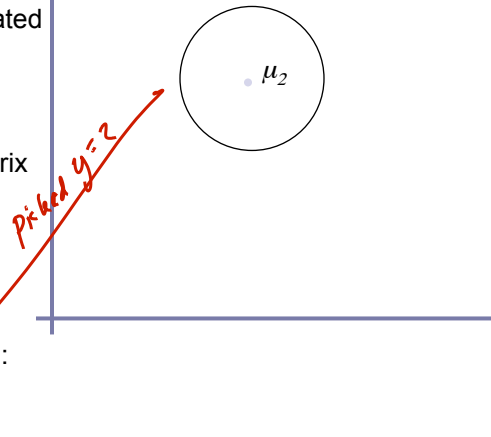
18

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 m_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 $P(y=i)$



©Carlos Guestrin 2005-2013

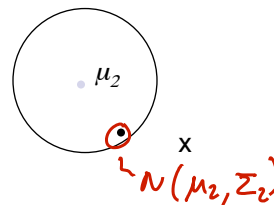
19

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 m_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 $\sim N(\mu_i, \sigma^2 I)$



©Carlos Guestrin 2005-2013

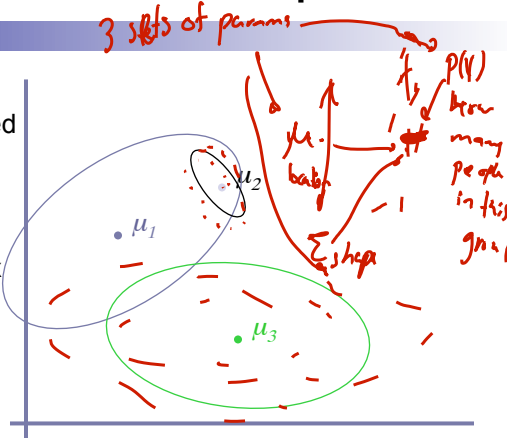
20

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



©Carlos Guestrin 2005-2013

21

Unsupervised Learning with Mixtures of Gaussians

EM Algorithm

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 17, 2013

©Carlos Guestrin 2005-2013

22

Supervised Learning of Mixtures of Gaussians

Mixtures of Gaussians:

□ Prior class probabilities: $P(y)$

□ Likelihood function per class: $P(x|y=i)$

$x \rightarrow m$ dims

■ Suppose, for each data point, we know location x and class y

□ Learning is easy... ☺

□ For prior $P(y)$

$$P(y=i) = \frac{\text{count}(y=i) \text{ in data}}{N}$$

□ For likelihood function:

$$P(x|y=i) = \frac{1}{\sigma_{ii}} \exp\left(-\frac{1}{2\sigma_{ii}} \sum_{j \in \text{cluster } i} (x_j^i - \mu_{ij})(x_j^i - \mu_{ij})\right)$$

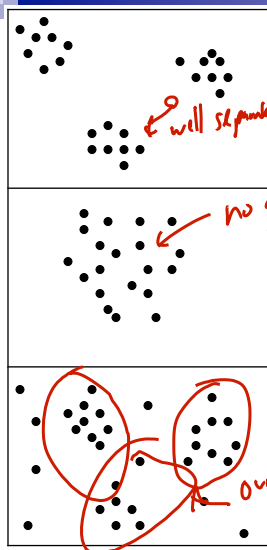
μ_i is average of x_j for points in class i

$\sigma_{ii} \leftarrow \sigma_{uv} = \frac{\sum_{j \in \text{cluster } i} (x_u^j - \mu_{iu})(x_v^j - \mu_{iv})}{\text{num points in cluster } i}$

©Carlos Guestrin 2005-2013

23

Unsupervised Learning: not as hard as it looks



Sometimes easy

Sometimes impossible

and sometimes in between

overlapping but meaningful clusters

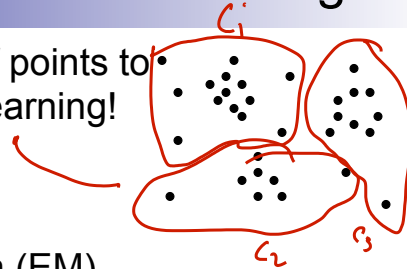
IN CASE YOU'RE WONDERING WHAT THESE DIAGRAMS ARE, THEY SHOW 2-d UNLABELED DATA (X VECTORS) DISTRIBUTED IN 2-d SPACE. THE TOP ONE HAS THREE VERY CLEAR GAUSSIAN CENTERS

©Carlos Guestrin 2005-2013

24

EM: “Reducing” Unsupervised Learning to Supervised Learning

- If we knew assignment of points to classes → Supervised Learning!



- Expectation-Maximization (EM)

- ☐ Guess assignment of points to classes *or clusters*
- ☐ Recompute model parameters
- ☐ Iterate

©Carlos Guestrin 2005-2013

25

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...

©Carlos Guestrin 2005-2013

26

Silly Example

Let events be "grades in a class"

$$\begin{array}{ll} w_1 = \text{Gets an A} & P(A) = \frac{1}{2} \\ w_2 = \text{Gets a B} & P(B) = \mu \\ w_3 = \text{Gets a C} & P(C) = 2\mu \\ w_4 = \text{Gets a D} & P(D) = \frac{1}{2} - 3\mu \end{array}$$

(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 ?

©Carlos Guestrin 2005-2013

27

Trivial Statistics

Normalization Constant

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

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

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

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

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

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!

©Carlos Guestrin 2005-2013

28

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 $1/2 : \mu$

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

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

MAXIMIZATION

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

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

REMEMBER

$$P(A) = 1/2$$

$$P(B) = \mu$$

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

$$P(D) = 1/2 - 3\mu$$

©Carlos Guestrin 2005-2013

29

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

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

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

E-step

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

M-step

= max like est. of μ given $b^{(t)}$

Continue iterating until ~~converged~~.

Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.

REMEMBER

$$P(A) = 1/2$$

$$P(B) = \mu$$

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

$$P(D) = 1/2 - 3\mu$$

©Carlos Guestrin 2005-2013

30

just like k-means

E.M. Convergence

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

Would we get the same answer with a diff. starting point?
Maybe... Maybe not...
Who knows,

In our example,
suppose we had

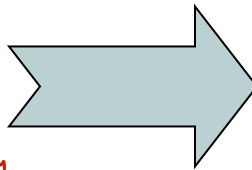
$$h = 20$$

$$c = 10$$

$$d = 10$$

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

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

20.317
got As