

EM

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

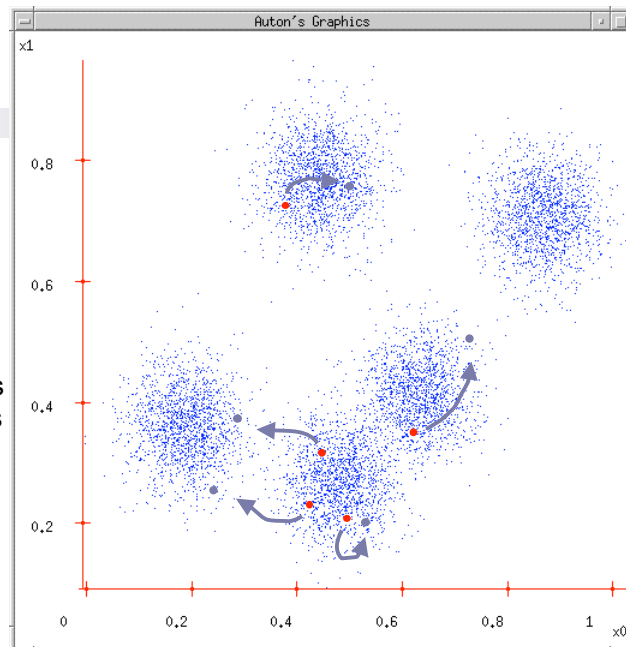
November 19th, 2007

©2005-2007 Carlos Guestrin

1

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!



2

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: *center of point j is closest to j*

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

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

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2$ *opt $\mu_i = \frac{\sum_{j: C^{(t)}(j)=i} x_j}{\sum_{j: C^{(t)}(j)=i} 1}$ is the mean!!*
 - Equivalent to $\mu_i \leftarrow$ average of its points!

Does K-means converge??? Part 2

- Optimize potential function:

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

- Fix C , optimize μ**

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

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

μ_i is mean of points in cluster i

\Rightarrow recenter in K-means.

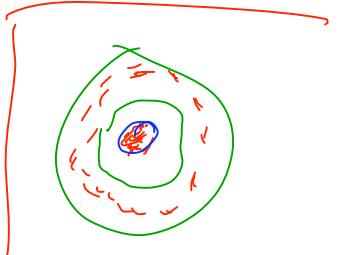
Coordinate descent algorithms

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

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

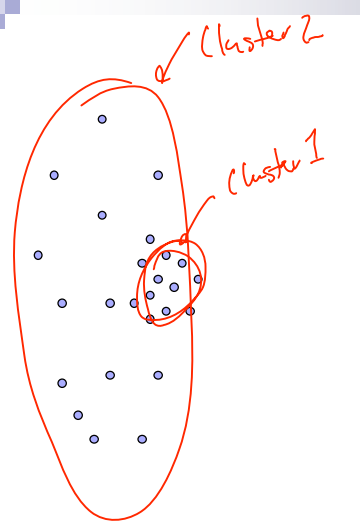
$F(a,b)$

$\nabla_{(a,b)} F(a,b) \approx$



5

(One) bad case for k-means



- Clusters may overlap
- Some clusters may be "wider" than others

6

Gaussian Bayes Classifier

Reminder

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

likelihood functions (pointing to $p(\mathbf{x}_j | y = i)$)
prior (pointing to $P(y = i)$)
multinomial (pointing to $P(y = i)$)

$P(Y=1) = 0.2$
 $P(Y=2) = 0.7$
 $P(Y=3) = 0.1$

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

constant (pointing to the fraction)
 $P(\mathbf{x}_j | y = i)$ (pointing to the exponential term)

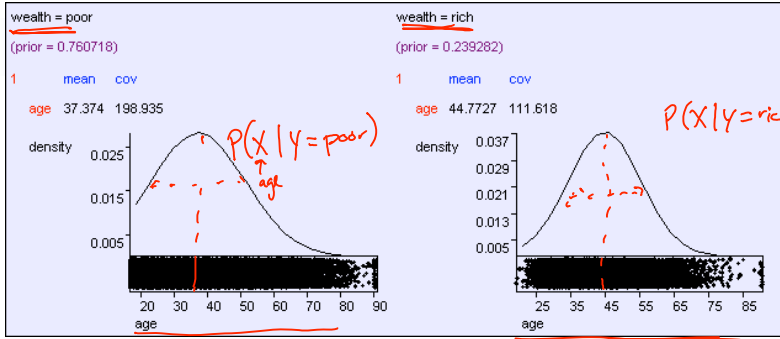
$$P(\mathbf{x}_j | y = i) \propto e^{-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i)}$$

$\mu_i = \begin{bmatrix} \mu_{i,1} \\ \mu_{i,2} \\ \vdots \\ \mu_{i,d} \end{bmatrix}$ *mean vector*
 $\mu_{i,2}$ *mean value of second dimension*

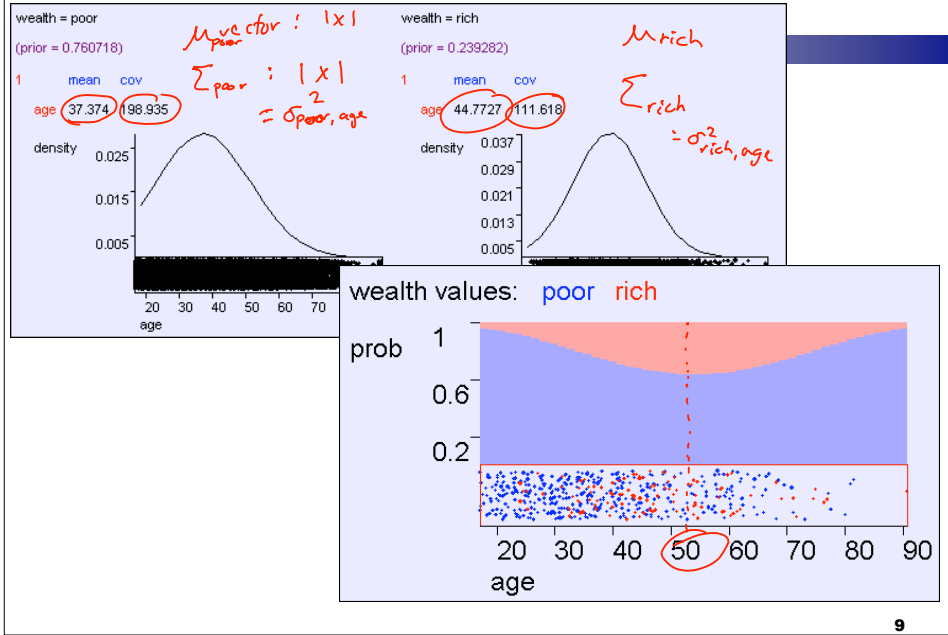
$\Sigma_i = \begin{bmatrix} \sigma_{i,11} & \sigma_{i,12} & \sigma_{i,13} \\ \sigma_{i,21} & \sigma_{i,22} & \dots \\ & & \ddots \\ & & & \sigma_{i,dd} \end{bmatrix}$ *covariance matrix*
covariance between features 1 and 2 (pointing to $\sigma_{i,12}$)

$\mathbf{x}_j \in \mathbb{R}^d$

Predicting wealth from age



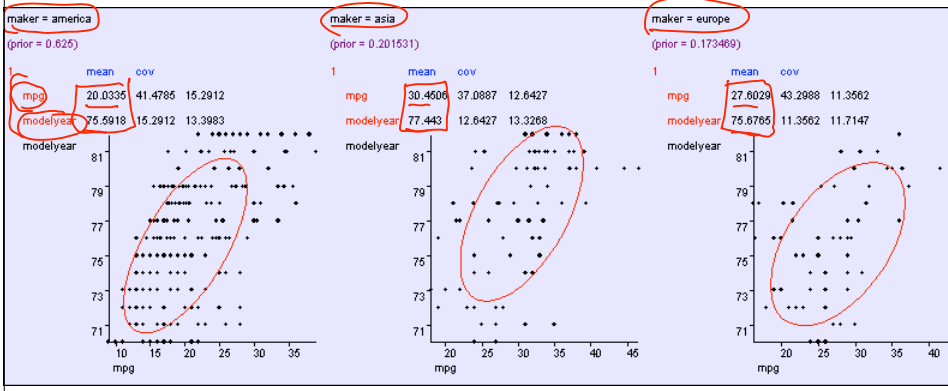
Predicting wealth from age



Learning modelyear, mpg ---> maker

$\Sigma_{\text{america}} = 2 \times 2$

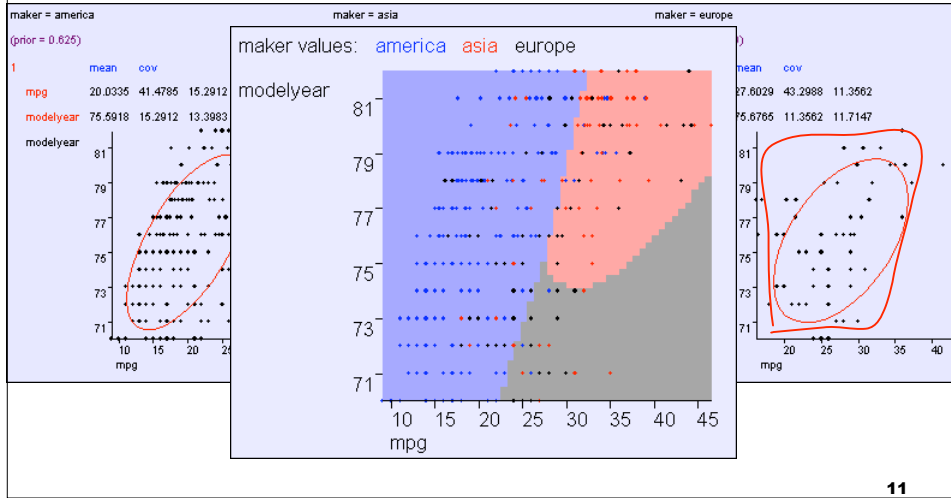
$$\Sigma_i = \begin{pmatrix} \sigma_{i,1}^2 & \sigma_{i,12} & \dots & \sigma_{i,1m} \\ \sigma_{i,12} & \sigma_{i,2}^2 & \dots & \sigma_{i,2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{i,1m} & \sigma_{i,2m} & \dots & \sigma_{i,m}^2 \end{pmatrix}$$



General: $O(m^2)$ parameters

per-class

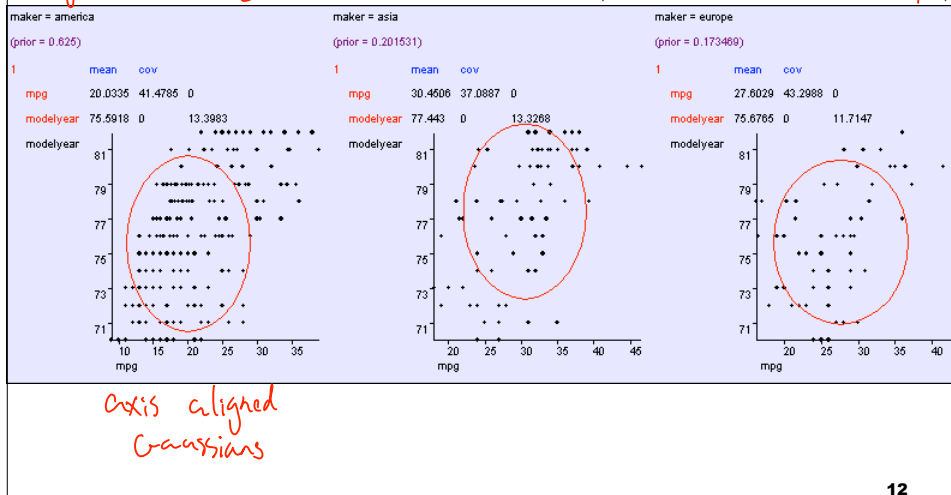
$$\Sigma_i = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{12} & \sigma_{22}^2 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1m} & \sigma_{2m} & \dots & \sigma_{mm}^2 \end{pmatrix}$$



Aligned: $O(m)$ parameters

Features are indep. given class [Gaussian NB]

$$\Sigma_i = \begin{pmatrix} \sigma_{i1}^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_{i2}^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_{i3}^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{im-1}^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_{im}^2 \end{pmatrix}$$

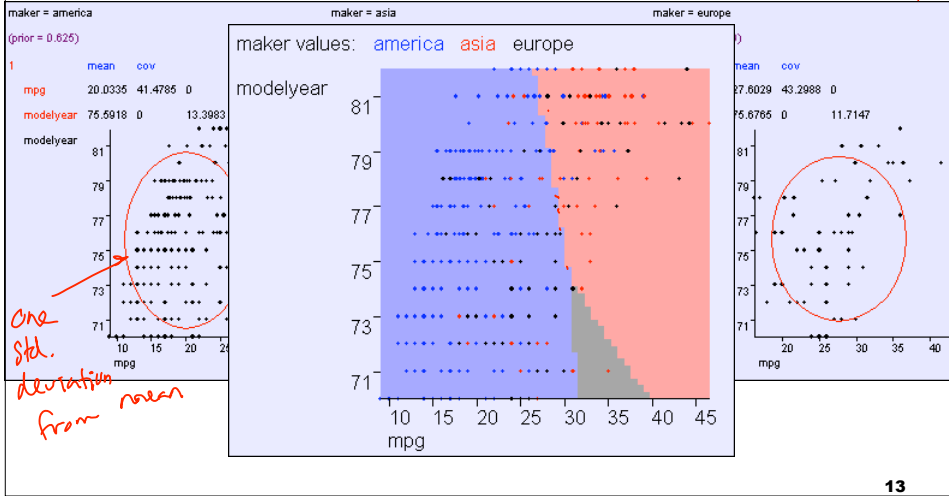


Aligned: $O(m)$
parameters

per class

$\Sigma_i =$

$$\begin{pmatrix} \sigma_{i1}^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_{i2}^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_{i3}^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_{im-1}^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_{im}^2 \end{pmatrix}$$

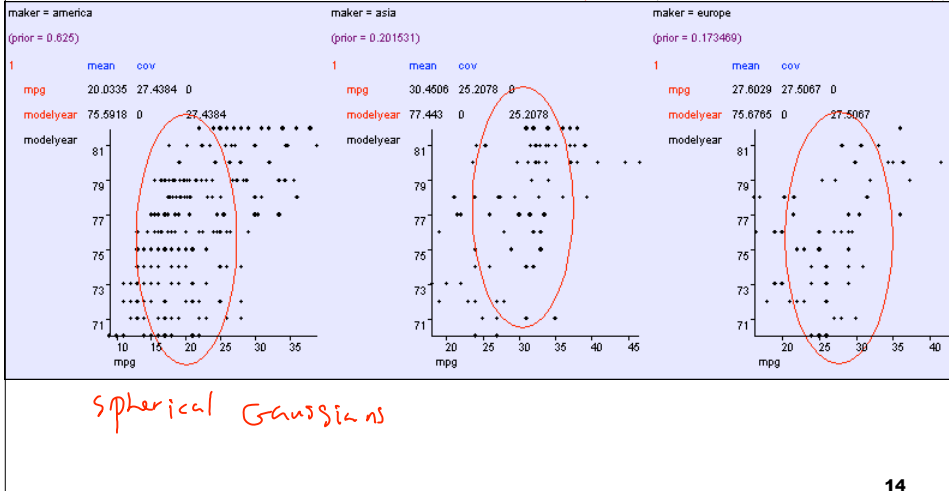


Spherical: $O(1)$
cov parameters

per class

every feature has same var.

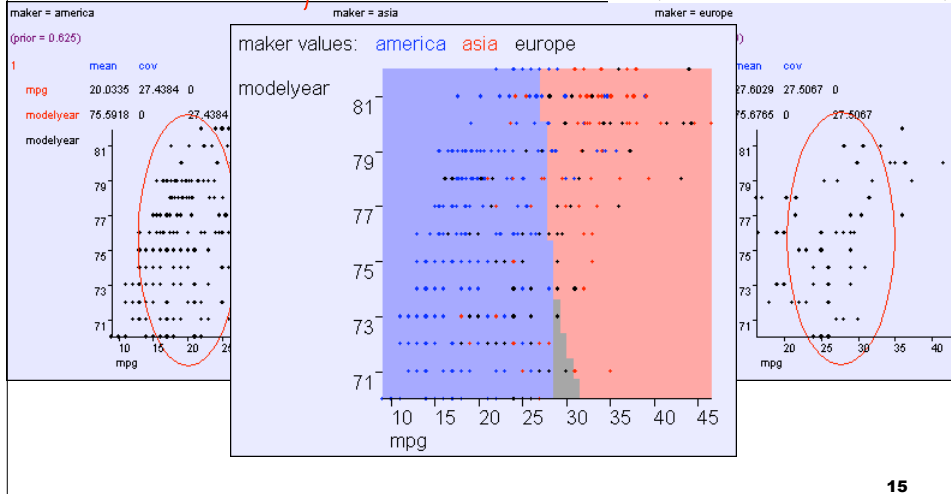
$$\Sigma = \begin{pmatrix} \sigma_i^2 & 0 & 0 & \dots & 0 & 0 \\ 0 & \sigma_i^2 & 0 & \dots & 0 & 0 \\ 0 & 0 & \sigma_i^2 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_i^2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \sigma_i^2 \end{pmatrix}$$



Spherical: $O(1)$
cov parameters

Carlos does not endorse, however....

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



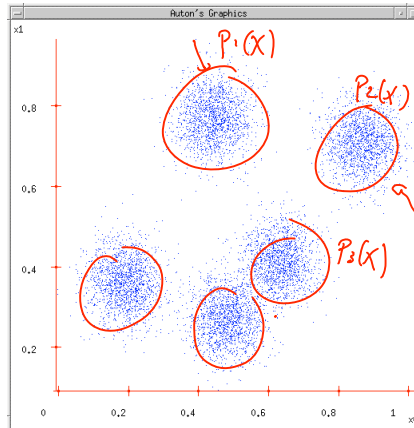
Next... back to Density Estimation

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

$$P(x) = \sum_i w_i P_i(x)$$

Gaussian mixture model

P_i is some gaussian with mean μ_i covariance Σ_i



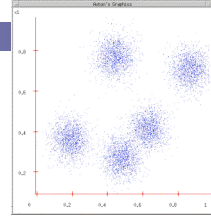
if data were labeled x_i, y_i for each i would know which group

But we don't see class labels!!!

- MLE:

- $\square \operatorname{argmax} \prod_j P(y_j, x_j)$

don't observe y



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

- Maximize marginal likelihood:

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

part I observe

Special case: spherical Gaussians and hard assignments

*$\sigma_{i,1}^2 = \sigma_{i,2}^2$
 $\Sigma_i = \begin{bmatrix} \sigma_{i,1}^2 & & 0 \\ & \sigma_{i,2}^2 & \\ 0 & & \sigma_{i,d}^2 \end{bmatrix}$
 $\Sigma_i = \Sigma_j$*

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

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

$P(x_j, y=i) = 0$ if $i \neq C(j)$

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

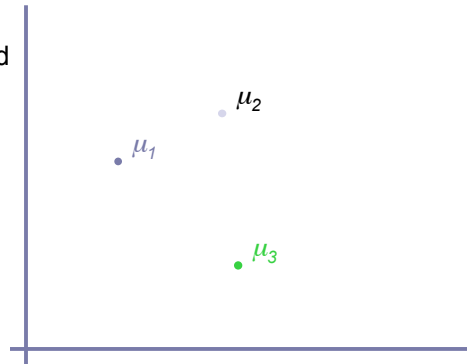
- Same as K-means!!!

$$= \sum_{j=1}^m \left[-\frac{1}{2\sigma^2} \|\mathbf{x}_j - \mu_{C(j)}\|^2\right] = -\frac{1}{2\sigma^2} \sum_{j=1}^m \|\mathbf{x}_j - \mu_{C(j)}\|^2$$

K-means objective

The GMM assumption

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

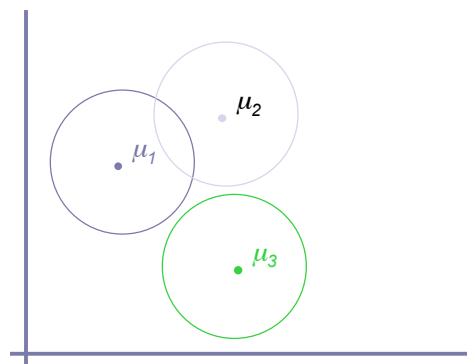


19

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 μ_i and covariance matrix $\sigma^2 I$

Each data point is generated according to the following recipe:



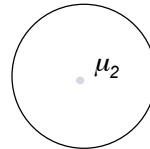
20

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 $P(y=i)$



21

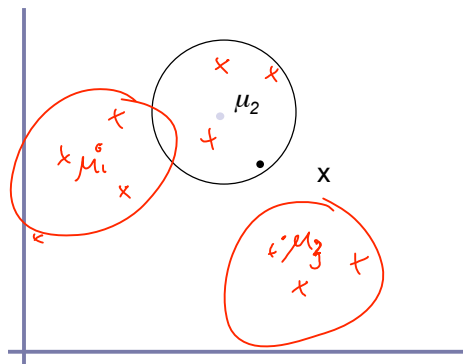
The GMM assumption

with
K-means
model

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



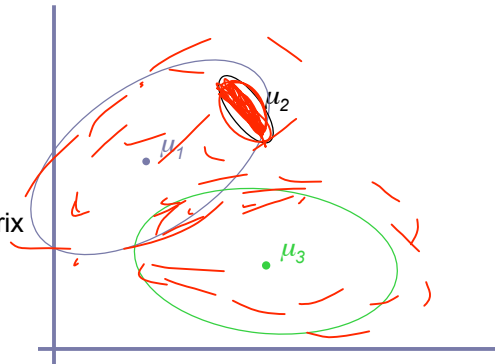
22

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

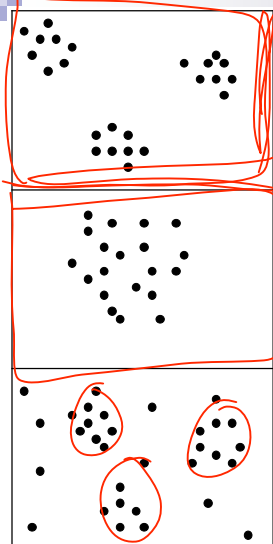
Each data point is generated according to the following recipe:

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



23

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

24

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: *sum over possible clusters*

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

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

likelihood *prior*

pick k using cross validation (generally will pick more clusters than you "want")

doesn't work for K-means ... !

Other criteria

maximize likelihood on held out data

25

Special case 2: spherical Gaussians and soft assignments

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

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

likelihood function

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

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

each point belongs to more than one cluster: prob. $P(y_j = i)$

26

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.

$P(y_j=i | x_j) \leftarrow$ Bayes rule

Suppose x 's are 1-dimensional.

(From Duda and Hart)

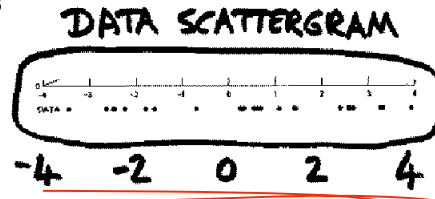
There are two classes; w_1 and w_2

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

Know these compute $P(x_j)$

There are 25 unlabeled datapoints

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

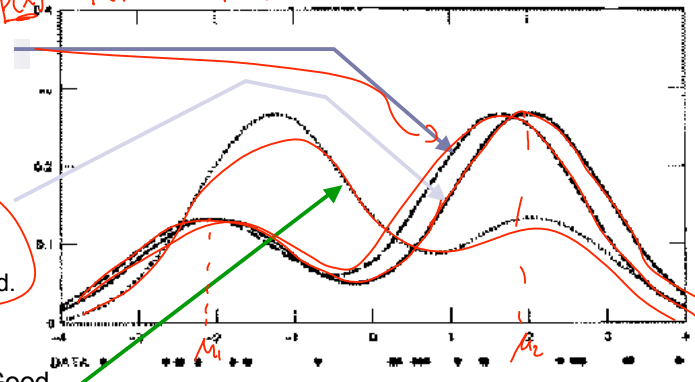


Duda & Hart's Example

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.

$P(x) = P(y_1) \cdot P(x|y=1) + P(y=2) \cdot P(x|y=2)$



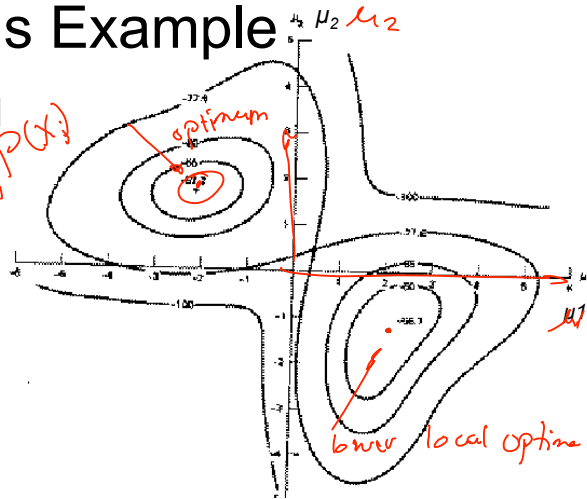
- They are close. Good.
- The 2nd solution tries to put the "2/3" hump where the "1/3" hump should go, and vice versa. *local optima*
- In this example unsupervised is almost as good as supervised. If the $x_1 \dots x_{25}$ are given the class which was used to learn them, then the results are $(\mu_1 = -2.176, \mu_2 = 1.684)$. Unsupervised got $(\mu_1 = -2.13, \mu_2 = 1.668)$.

Duda & Hart's Example



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

*max marginal (i.e. $P(x_i)$)
 μ_1, μ_2 likelihood*



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 $P(\text{data} | \mu_1, \mu_2 \dots \mu_k)$

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

- The normal max likelihood trick:
 Set $\frac{\partial}{\partial \mu_i} \log \text{Prob}(\dots) = 0$
 and solve for μ_i 's.
 # Here you get non-linear non-analytically-solvable equations
- Use gradient descent
 Often slow but doable
- Use a much faster, cuter, and recently very popular method...

*not convex
 (previous slide)*

Announcements

■ HW5 out later today...

- Due December 5th by 3pm to Monica Hopes, Wean 4619

■ Project:

- Poster session: NSH Atrium, Friday 11/30, 2-5pm

■ Print your poster early!!!

- SCS facilities has a poster printer, ask helpdesk
- Students from outside SCS should check with their departments
- It's OK to print separate pages

■ We'll provide pins, posterboard and an easel

- Poster size: 32x40 inches

■ Invite your friends, there will be a prize for best poster, by popular vote

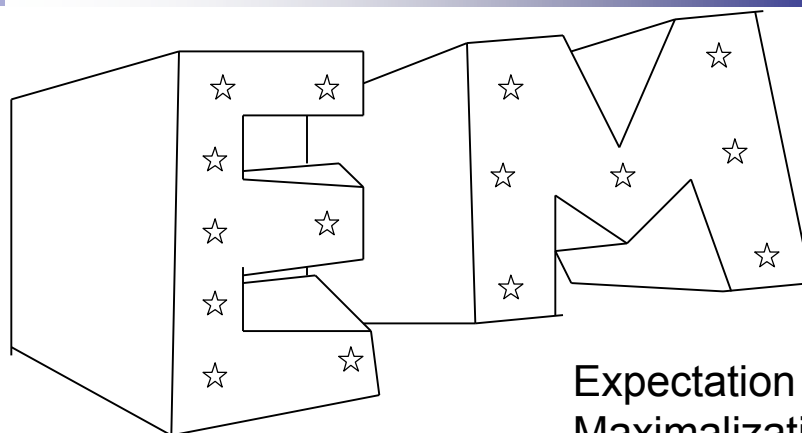
■ Last lecture:

- Thursday, 11/29, 5-6:20pm, Wean 7500

4:45pm

*arrive
15 mins
early to
set up!*

31



32

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 learning BNs with hidden data

33

Silly Example

Let events be "grades in a class"

$w_1 =$ Gets an <u>A</u>	$P(A) = \frac{1}{2}$
$w_2 =$ Gets a <u>B</u>	$P(B) = \mu$
$w_3 =$ Gets a <u>C</u>	$P(C) = 2\mu$
$w_4 =$ Gets a <u>D</u>	$P(D) = \frac{1}{2} - 3\mu$

(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 ?

$\hat{\mu}_{MLE}$

34

Trivial Statistics

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

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

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

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!

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?

A = \bar{a} people

B = $h - \bar{a}$ people

what's ?

if I knew \bar{a}

$\mu = \frac{h - \bar{a} + c}{6(h - \bar{a} + c + d)}$

What's \bar{a} ?

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

iterate

Since the ratio a:b should be the same as the ratio $\frac{1}{2} : \mu$

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

MAXIMIZATION

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

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

REMEMBER

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

P(B) = μ

P(C) = 2μ

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

37

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

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

if true param was $\mu^{(t)}$

E-step

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

= max like est. of μ given $b^{(t)}$

M-step

Continue iterating until converged.

Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.

REMEMBER

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

P(B) = μ

P(C) = 2μ

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

38

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]

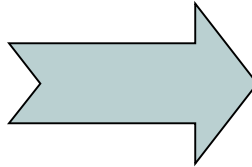
In our example,
suppose we had

$$h = 20$$

$$c = 10$$

$$d = 10$$

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



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

$a(t)$
20