

So far, supervised learning:  
 $h: \mathcal{X} \rightarrow \mathcal{R}$  ← regression  
 $h: \mathcal{X} \rightarrow \{0, 1, \dots, k\}$  ← classification

Unsupervised learning

# Clustering

## K-means Continued

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 15, 2013

©Carlos Guestrin 2005-2013

1

## Clustering images

no labels given



beach



flowers



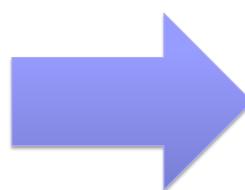
sky

;

;

;

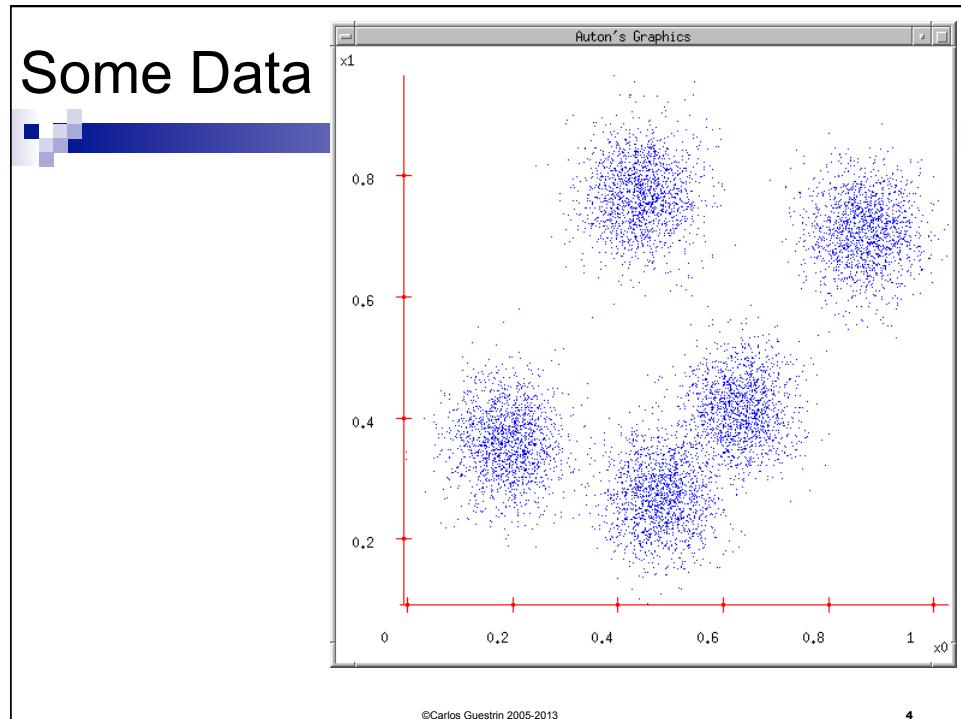
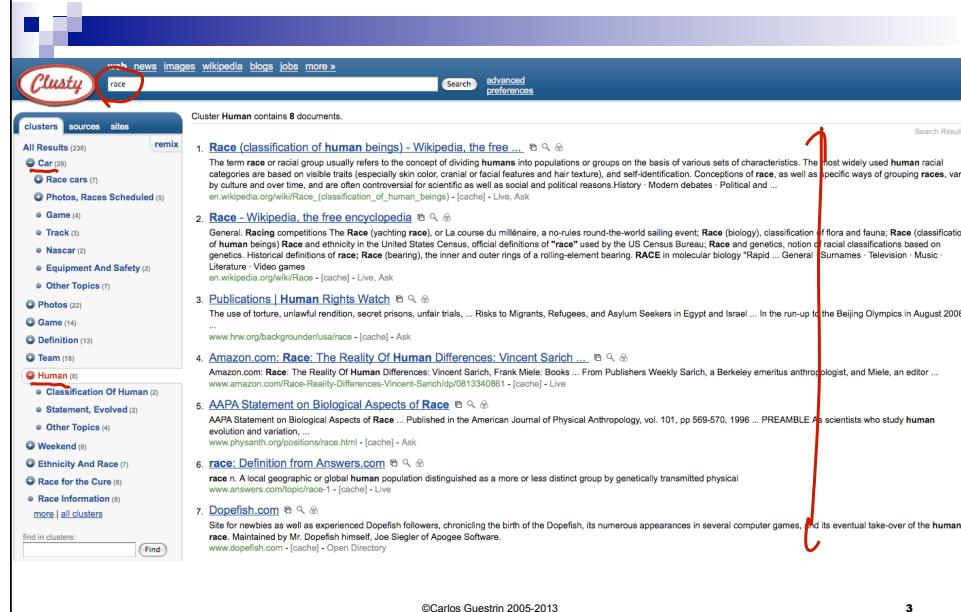
Set of Images



©Carlos Guestrin 2005-2013

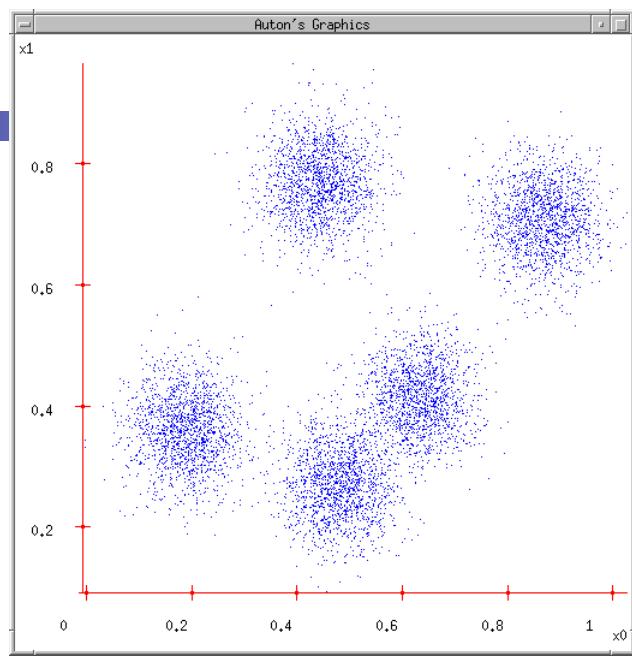
[Goldberger et al.] 2

# Clustering web search results



# K-means

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

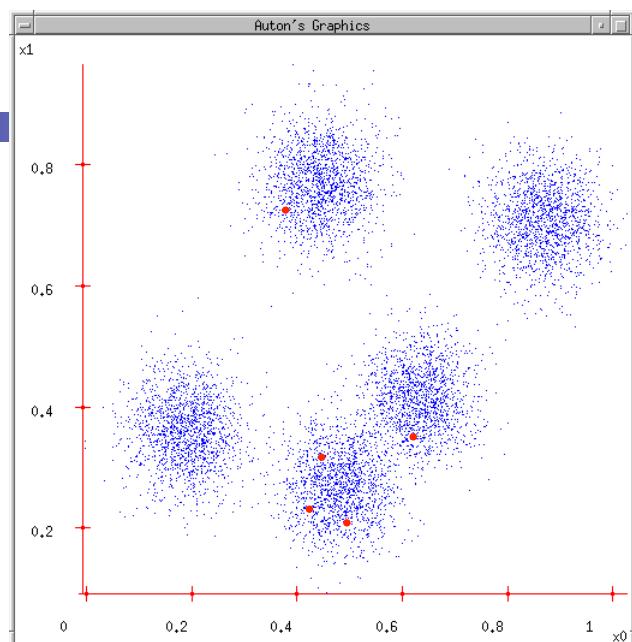


©Carlos Guestrin 2005-2013

5

# K-means

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

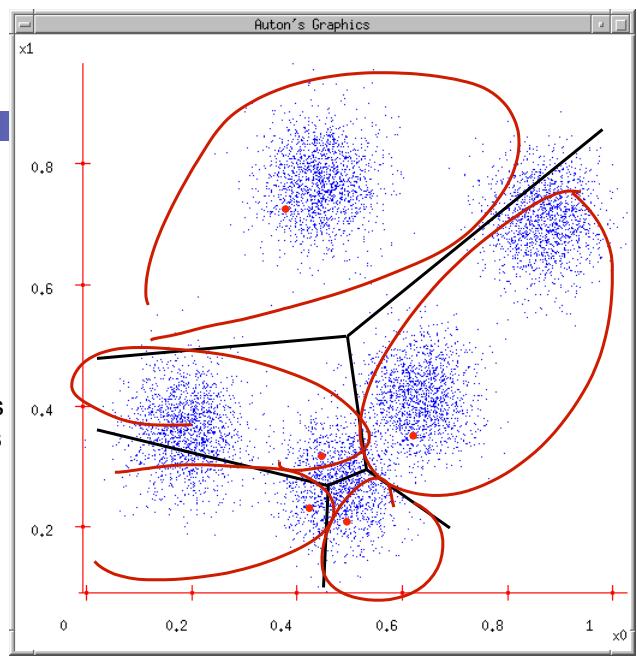


©Carlos Guestrin 2005-2013

6

# 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. (Thus each Center "owns" a set of datapoints)

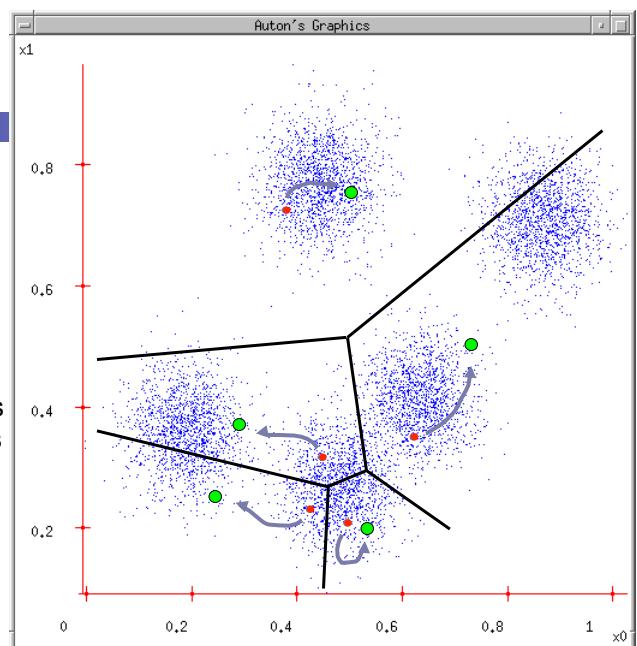


©Carlos Guestrin 2005-2013

7

# 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



©Carlos Guestrin 2005-2013

8

## 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!

©Carlos Guestrin 2005-2013

9

## K-means

- Randomly initialize  $k$  centers *or smartly*
  - $\square \mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$
  - Converge when nothing moves (no point changes its cluster)*
- Classify: Assign each point  $j \in \{1, \dots, n\}$  to nearest center:
  - $\square C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$  *in cluster center* *data point* *fix  $\mu$ , OPTC*
- Recenter:  $\mu_i$  becomes centroid of its point:
  - $\square \mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j:C(j)=i} \|\mu - x_j\|^2 \leftarrow \mu_i = \frac{\sum_{j:C(j)=i} x_j}{\text{num of points assigned to cluster } i}$
  - $\square$  Equivalent to  $\mu_i \leftarrow \text{average of its points!}$

©Carlos Guestrin 2005-2013

10

## What is K-means optimizing?

- Potential function  $F(\mu, C)$  of centers  $\mu$  and point allocations  $C$ : *distance over all points*

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

↑ *k centers*  
 ↑ *allocations*

- Optimal K-means:

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

*coordinate descent*  
*optimization is hard*

©Carlos Guestrin 2005-2013

11

## 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  $C$

$$\begin{aligned}
 & \min_{C: C(1) \dots C(N)} \sum_{j=1}^N \|\mu_{C(j)} - x_j\|^2 = \min_{C(1)} \min_{C(2)} \dots \min_{C(k)} \sum_{j=1}^N \|\mu_{C(j)} - x_j\|^2 \\
 & = \sum_{j=1}^N \min_{C(j)} \|\mu_{C(j)} - x_j\|^2
 \end{aligned}$$

*indep.*  
*minimization*      *exactly*      *the "classification step"*

©Carlos Guestrin 2005-2013

12

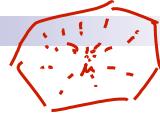
## Does K-means converge??? Part 2

- 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  $\mu$

$$\begin{aligned} \min_{\mu: \mu_1, \dots, \mu_K} \sum_{j=1}^N \|\mu_{c(j)} - x_j\|^2 &= \min_{\mu: \mu_1, \dots, \mu_K} \sum_{i=1}^k \sum_{j:c(j)=i} \|\mu_i - x_j\|^2 \\ &\equiv \sum_{i=1}^k \min_{\mu_i} \sum_{j:c(j)=i} \|\mu_i - x_j\|^2 \\ \mu_i &= \text{center of points} = \text{average} = \frac{\sum_{j:c(j)=i} x_j}{\text{number of points in cluster } i} \end{aligned}$$



©Carlos Guestrin 2005-2013

13

## Coordinate descent algorithms

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

$$\begin{aligned} \min_{\mu} \min_C F(\mu, C) &= \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2 \\ a, b &\cdot F(a, b) \\ a' &\downarrow \min \\ a', b' &\quad F(a', b') \leq F(a, b) \\ a', b' &\downarrow \min \\ a', b' &\quad F(a', b') \leq F(a', b) \leq F(a, b) \end{aligned}$$

$\geq 0$   
non-increasing sequence  
lower bound

- Coordinate descent:

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

- Converges!!!

- if  $F$  is bounded below
- to a (often good) local optimum
  - as we saw in applet (play with it!)
  - (For LASSO it converged to the optimum)

Random restarts help

- K-means is a coordinate descent algorithm!

©Carlos Guestrin 2005-2013

14

# Mixtures of Gaussians

Machine Learning – CSE446

Carlos Guestrin

University of Washington

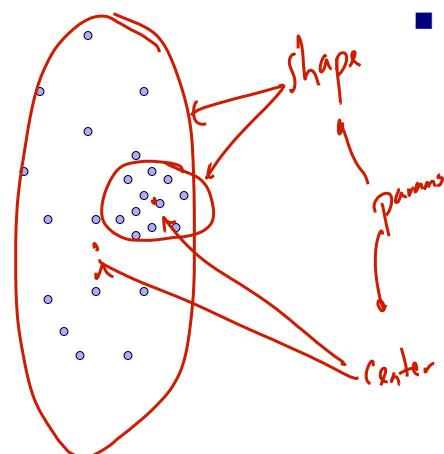
May 15, 2013

©Carlos Guestrin 2005-2013

15

## (One) bad case for k-means

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

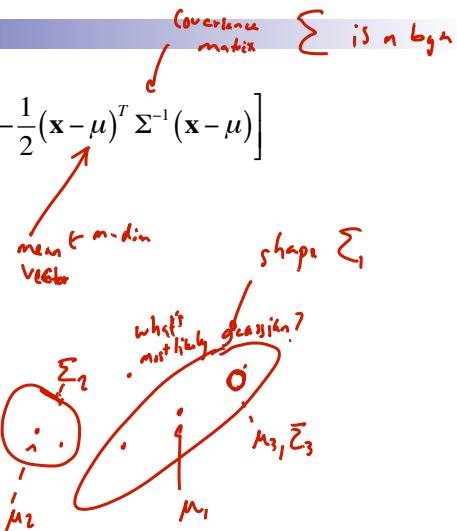
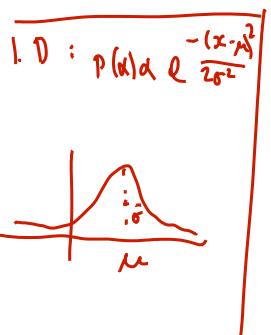


©Carlos Guestrin 2005-2013

16

## Gaussians in $m$ Dimensions

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



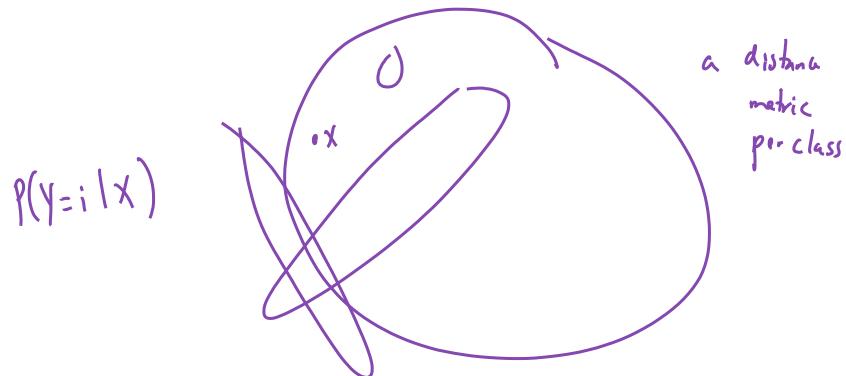
©Carlos Guestrin 2005-2013

17

Suppose You Have a Gaussian For Each Class

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

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



©Carlos Guestrin 2005-2013

18

# Gaussian Bayes Classifier

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

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

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

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

- Thank you Bayes Rule!!

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

$$\propto \frac{1}{(2\pi)^{m/2} |\Sigma_i|^{1/2}} \exp^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)} \cdot P(y = i)$$

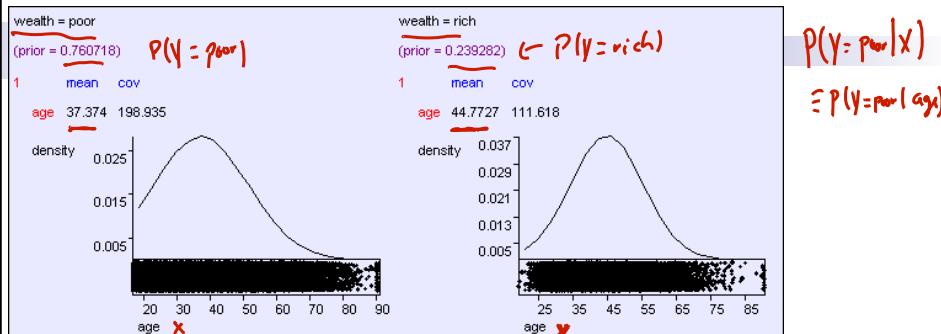
(cluster data likelihoods)

prior: before we see  $\mathbf{x}$ , what fraction of points we expect to fall in cluster  $i$

©Carlos Guestrin 2005-2013

19

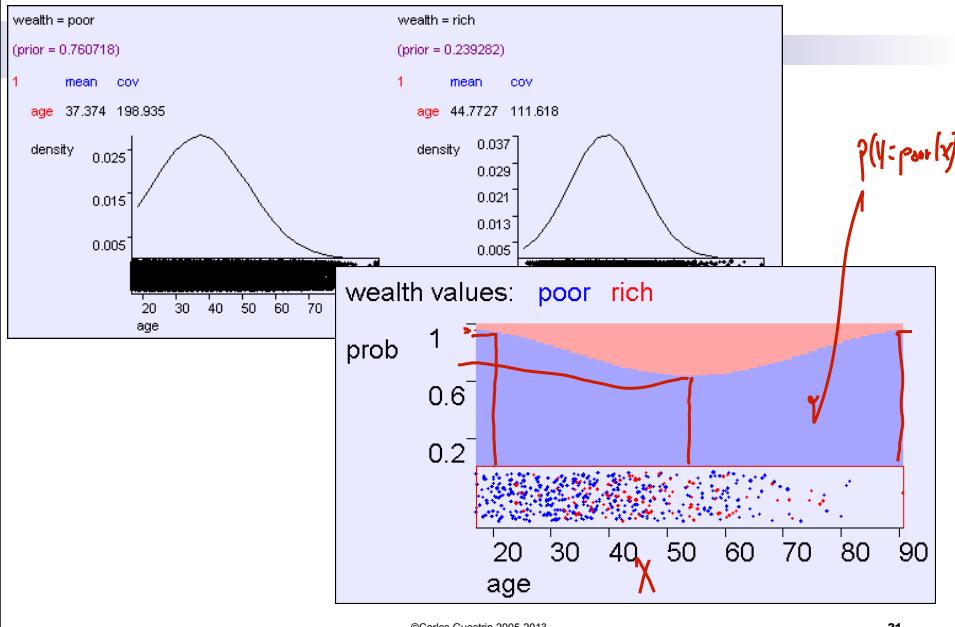
## Predicting wealth from age



©Carlos Guestrin 2005-2013

20

# Predicting wealth from age

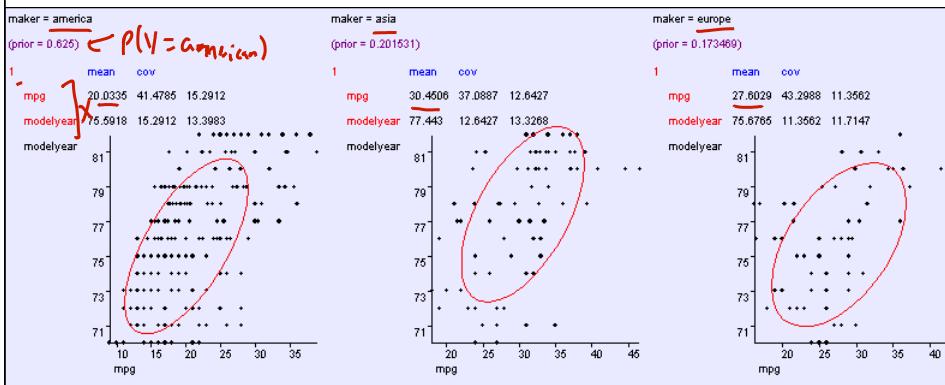


©Carlos Guestrin 2005-2013

21

## Learning modelyear , mpg ---> maker

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



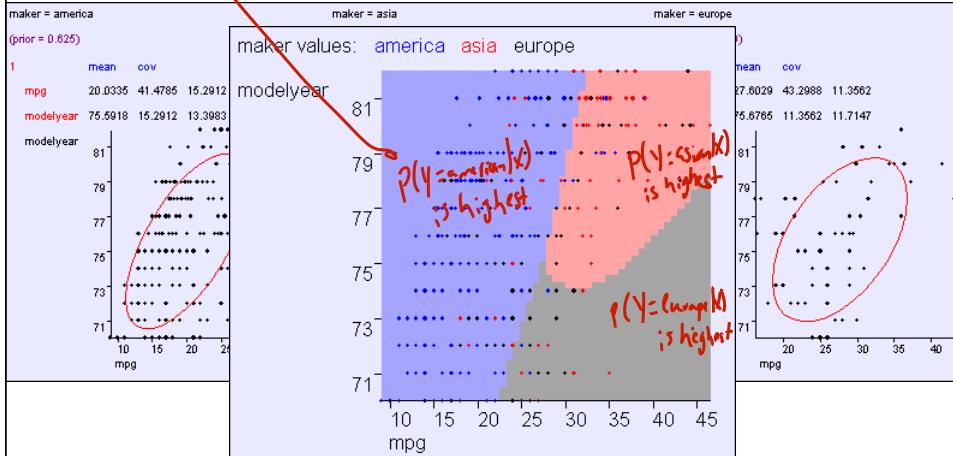
©Carlos Guestrin 2005-2013

22

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

*(classification  $\arg\max_y P(y|x)$ )*

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



©Carlos Guestrin 2005-2013

23