

# Probability Densities in Data Mining

**Andrew W. Moore**  
**Associate Professor**  
**School of Computer Science**  
**Carnegie Mellon University**

[www.cs.cmu.edu/~awm](http://www.cs.cmu.edu/~awm)  
[awm@cs.cmu.edu](mailto:awm@cs.cmu.edu)  
412-268-7599

Copyright © 2001, Andrew W. Moore

Aug 27, 2001

## Probability Densities in Data Mining

- Why we should care
- Notation and Fundamentals of continuous PDFs
- Multivariate continuous PDFs
- Combining continuous and discrete random variables

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 2

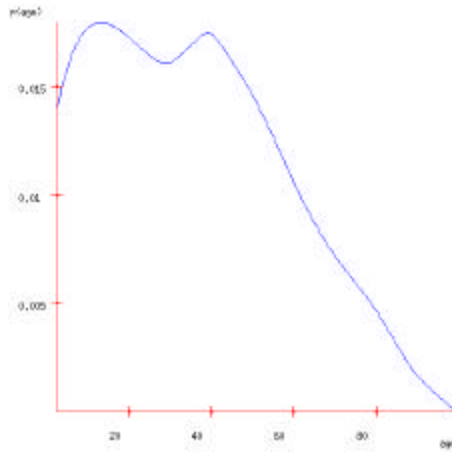
## Why we should care

- Real Numbers occur in at least 50% of database records
- Can't always quantize them
- So need to understand how to describe where they come from
- A great way of saying what's a reasonable range of values
- A great way of saying how multiple attributes should reasonably co-occur

## Why we should care

- Can immediately get us Bayes Classifiers that are sensible with real-valued data
- You'll need to **intimately** understand PDFs in order to do kernel methods, clustering with Mixture Models, analysis of variance, time series and many other things
- Will introduce us to linear and non-linear regression

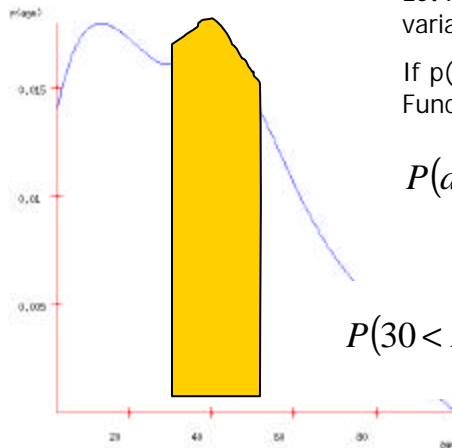
## A PDF of American Ages in 2000



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 5

## A PDF of American Ages in 2000



Let  $X$  be a continuous random variable.

If  $p(x)$  is a Probability Density Function for  $X$  then...

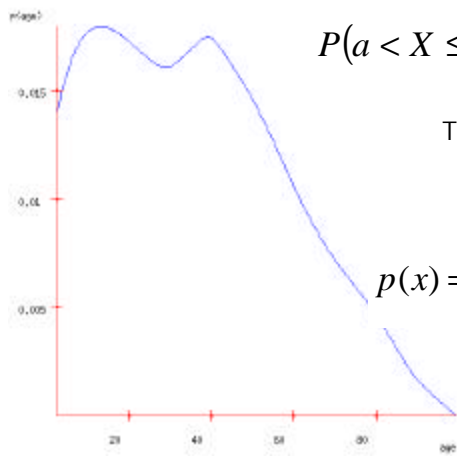
$$P(a < X \leq b) = \int_{x=a}^b p(x) dx$$

$$P(30 < \text{Age} \leq 50) = \int_{\text{age}=30}^{50} p(\text{age}) d\text{age} = 0.36$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 6

# Properties of PDFs



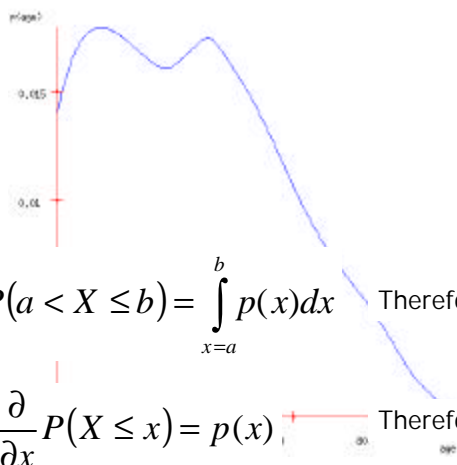
$$P(a < X \leq b) = \int_{x=a}^b p(x) dx$$

That means...

$$p(x) = \lim_{h \rightarrow 0} \frac{P\left(x - \frac{h}{2} < X \leq x + \frac{h}{2}\right)}{h}$$

$$\frac{\partial}{\partial x} P(X \leq x) = p(x)$$

# Properties of PDFs



$$P(a < X \leq b) = \int_{x=a}^b p(x) dx$$

Therefore...

$$\int_{x=-\infty}^{\infty} p(x) dx = 1$$

$$\frac{\partial}{\partial x} P(X \leq x) = p(x)$$

Therefore...

$$\forall x : p(x) \geq 0$$

## Talking to your stomach

- What's the gut-feel meaning of  $p(x)$ ?

If

$$p(5.31) = 0.06 \text{ and } p(5.92) = 0.03$$

then

when a value  $X$  is sampled from the distribution, you are 2 times as likely to find that  $X$  is "very close to" 5.31 than that  $X$  is "very close to" 5.92.

## Talking to your stomach

- What's the gut-feel meaning of  $p(x)$ ?

If

$$p(a) = 0.06 \text{ and } p(b) = 0.03$$

then

when a value  $X$  is sampled from the distribution, you are 2 times as likely to find that  $X$  is "very close to"  $a$  than that  $X$  is "very close to"  $b$ .

## Talking to your stomach

- What's the gut-feel meaning of  $p(x)$ ?

If

$$p(a) = 2z \quad \text{and} \quad p(b) = z$$

then

when a value  $X$  is sampled from the distribution, you are 2 times as likely to find that  $X$  is "very close to"  $a$  than that  $X$  is "very close to"  $b$ .

## Talking to your stomach

- What's the gut-feel meaning of  $p(x)$ ?

If

$$p(a) = az \quad \text{and} \quad p(b) = z$$

then

when a value  $X$  is sampled from the distribution, you are  $a$  times as likely to find that  $X$  is "very close to"  $a$  than that  $X$  is "very close to"  $b$ .

## Talking to your stomach

- What's the gut-feel meaning of  $p(x)$ ?

If 
$$\frac{p(a)}{p(b)} = a$$

then

when a value  $X$  is sampled from the distribution, you are  $a$  times as likely to find that  $X$  is "very close to"  $a$  than that  $X$  is "very close to"  $b$ .

## Talking to your stomach

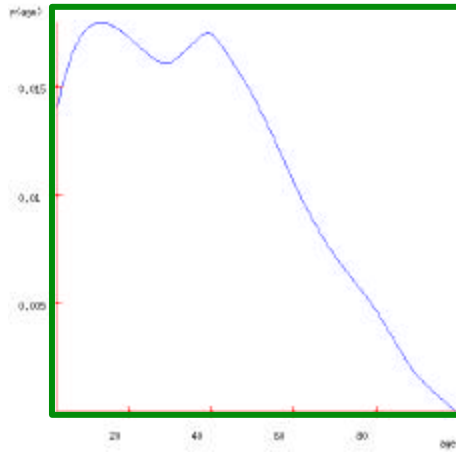
- What's the gut-feel meaning of  $p(x)$ ?

If 
$$\frac{p(a)}{p(b)} = a$$

then

$$\lim_{h \rightarrow 0} \frac{P(a-h < X < a+h)}{P(b-h < X < b+h)} = a$$

## Yet another way to view a PDF



A recipe for sampling a random age.

1. Generate a random dot from the rectangle surrounding the PDF curve. Call the dot (age,d)
2. If  $d < p(\text{age})$  stop and return age
3. Else try again: go to Step 1.

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 15

## Test your understanding

- True or False:

$$\forall x: p(x) \leq 1$$

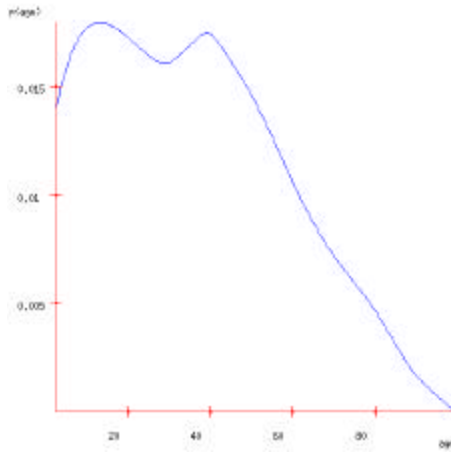
$$\forall x: P(X = x) = 0$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 16



# Expectations



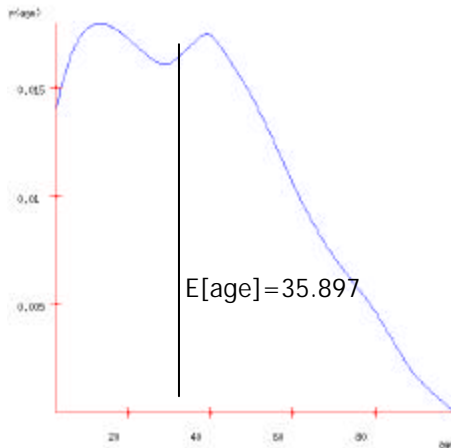
$E[X]$  = the expected value of random variable  $X$   
= the average value we'd see if we took a very large number of random samples of  $X$

$$= \int_{x=-\infty}^{\infty} x p(x) dx$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 17

# Expectations



$E[X]$  = the expected value of random variable  $X$   
= the average value we'd see if we took a very large number of random samples of  $X$

$$= \int_{x=-\infty}^{\infty} x p(x) dx$$

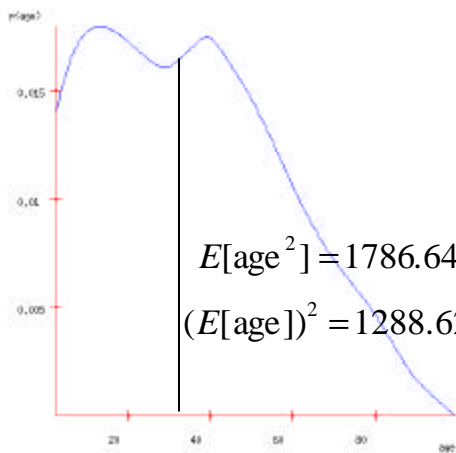
= the first moment of the shape formed by the axes and the blue curve

= the best value to choose if you must guess an unknown person's age and you'll be fined the square of your error

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 18

## Expectation of a function



$\mu = E[f(X)]$  = the expected value of  $f(x)$  where  $x$  is drawn from  $X$ 's distribution.

= the average value we'd see if we took a very large number of random samples of  $f(X)$

$$m = \int_{x=-\infty}^{\infty} f(x) p(x) dx$$

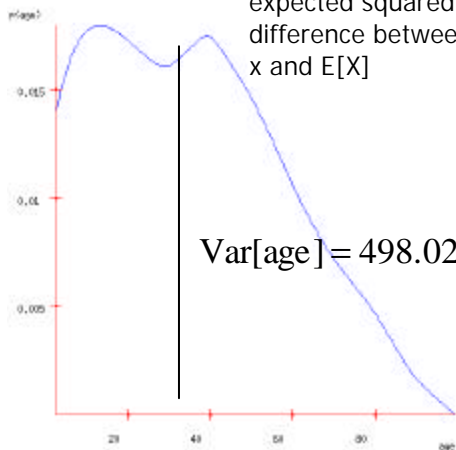
Note that in general:

$$E[f(x)] \neq f(E[X])$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 19

## Variance



$\sigma^2 = \text{Var}[X]$  = the expected squared difference between  $x$  and  $E[X]$

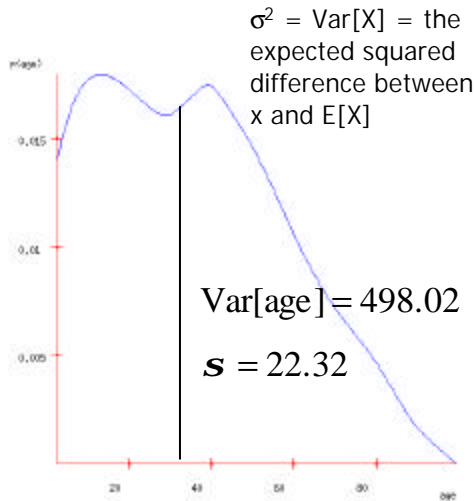
$$s^2 = \int_{x=-\infty}^{\infty} (x - m)^2 p(x) dx$$

= amount you'd expect to lose if you must guess an unknown person's age and you'll be fined the square of your error, and assuming you play optimally

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 20

# Standard Deviation



$\sigma^2 = \text{Var}[X]$  = the expected squared difference between  $x$  and  $E[X]$

$$s^2 = \int_{-\infty}^{\infty} (x - m)^2 p(x) dx$$

= amount you'd expect to lose if you must guess an unknown person's age and you'll be fined the square of your error, and assuming you play optimally

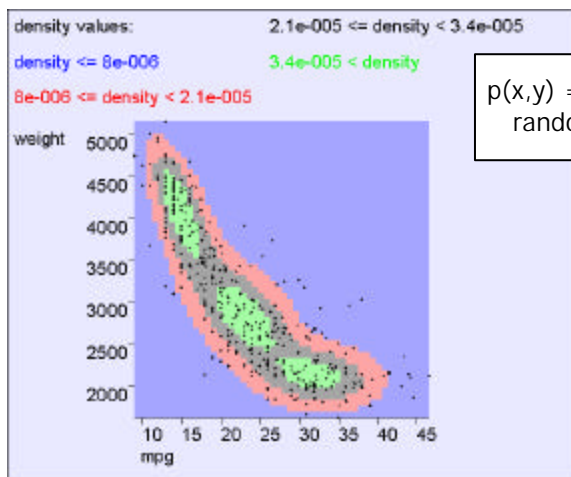
$\sigma$  = Standard Deviation = "typical" deviation of  $X$  from its mean

$$s = \sqrt{\text{Var}[X]}$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 21

# In 2 dimensions



$p(x,y)$  = probability density of random variables  $(X,Y)$  at location  $(x,y)$

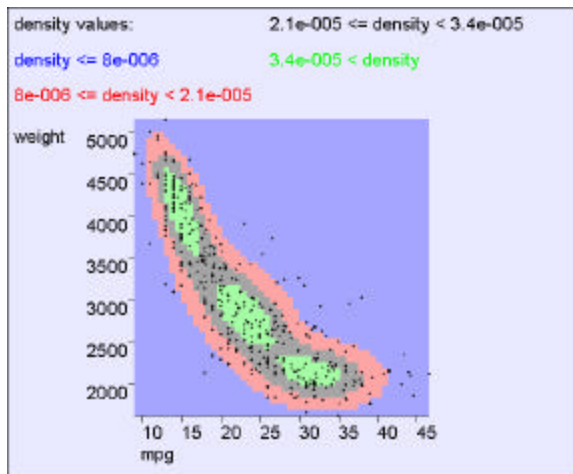
Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 22

## In 2 dimensions

Let  $X, Y$  be a pair of continuous random variables, and let  $R$  be some region of  $(X, Y)$  space...

$$P((X, Y) \in R) = \iint_{(x, y) \in R} p(x, y) dy dx$$



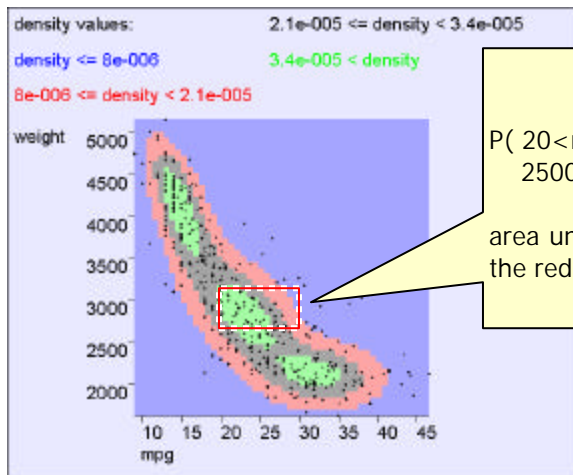
Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 23

## In 2 dimensions

Let  $X, Y$  be a pair of continuous random variables, and let  $R$  be some region of  $(X, Y)$  space...

$$P((X, Y) \in R) = \iint_{(x, y) \in R} p(x, y) dy dx$$



Copyright © 2001, Andrew W. Moore

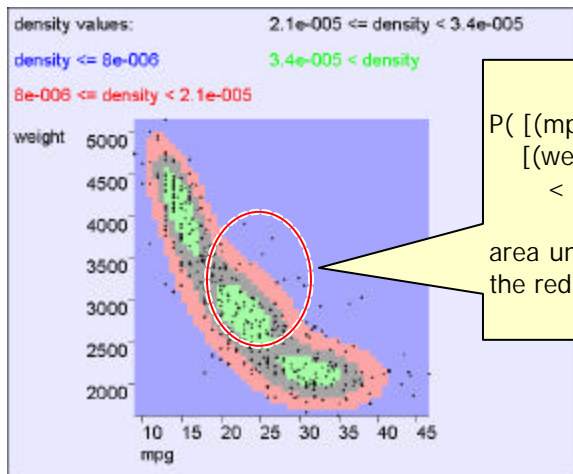
Probability Densities: Slide 24

$P(20 < \text{mpg} < 30 \text{ and } 2500 < \text{weight} < 3000) =$   
area under the 2-d surface within the red rectangle

## In 2 dimensions

Let  $X, Y$  be a pair of continuous random variables, and let  $R$  be some region of  $(X, Y)$  space...

$$P((X, Y) \in R) = \iint_{(x, y) \in R} p(x, y) dy dx$$



$$P\left(\left[\frac{\text{mpg}-25}{10}\right]^2 + \left[\frac{\text{weight}-3300}{1500}\right]^2 < 1\right) =$$

area under the 2-d surface within the red oval

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 25

## In 2 dimensions

Let  $X, Y$  be a pair of continuous random variables, and let  $R$  be some region of  $(X, Y)$  space...

$$P((X, Y) \in R) = \iint_{(x, y) \in R} p(x, y) dy dx$$

Take the special case of region  $R = \text{"everywhere"}$ .

Remember that with probability 1,  $(X, Y)$  will be drawn from "somewhere".

So..

$$\int_{x=-\infty}^{\infty} \int_{y=-\infty}^{\infty} p(x, y) dy dx = 1$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 26

## In 2 dimensions

Let  $X, Y$  be a pair of continuous random variables, and let  $R$  be some region of  $(X, Y)$  space...

$$P((X, Y) \in R) = \iint_{(x, y) \in R} p(x, y) dy dx$$

$$p(x, y) = \lim_{h \rightarrow 0} \frac{P\left(x - \frac{h}{2} < X \leq x + \frac{h}{2} \quad \wedge \quad y - \frac{h}{2} < Y \leq y + \frac{h}{2}\right)}{h^2}$$

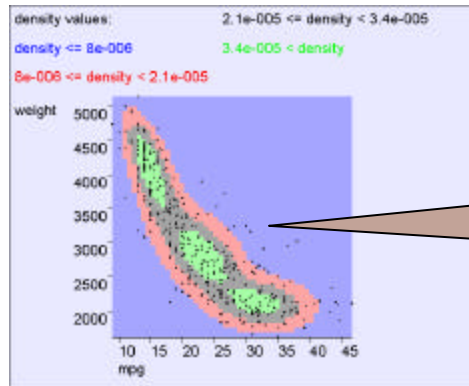
## In $m$ dimensions

Let  $(X_1, X_2, \dots, X_m)$  be an  $n$ -tuple of continuous random variables, and let  $R$  be some region of  $\mathbf{R}^m$  ...

$$P((X_1, X_2, \dots, X_m) \in R) = \iiint_{(x_1, x_2, \dots, x_m) \in R} \dots \int p(x_1, x_2, \dots, x_m) dx_m, \dots, dx_2, dx_1$$

# Independence

$$X \perp Y \text{ iff } \forall x, y: p(x, y) = p(x)p(y)$$



If X and Y are independent then knowing the value of X does not help predict the value of Y

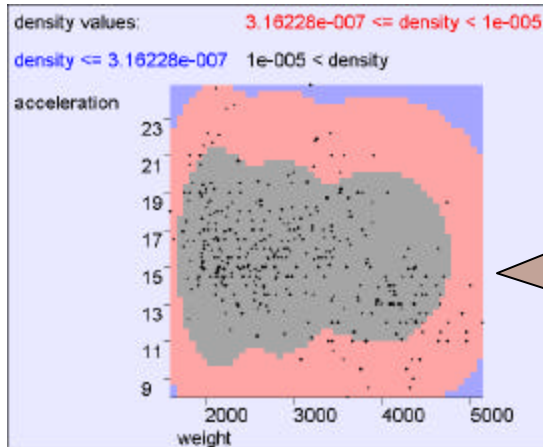
mpg, weight NOT independent

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 29

# Independence

$$X \perp Y \text{ iff } \forall x, y: p(x, y) = p(x)p(y)$$



If X and Y are independent then knowing the value of X does not help predict the value of Y

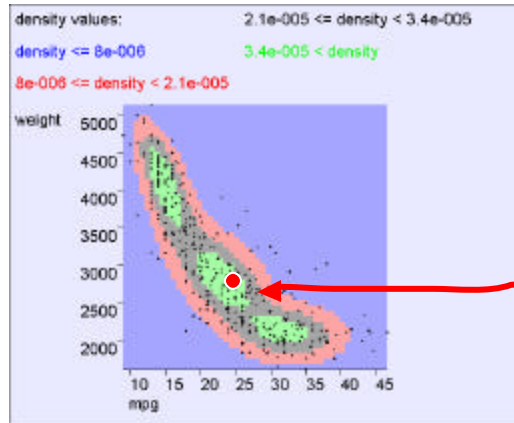
the contours say that acceleration and weight are independent

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 30

## Multivariate Expectation

$$\boldsymbol{\mu}_{\mathbf{X}} = E[\mathbf{X}] = \int \mathbf{x} p(\mathbf{x}) d\mathbf{x}$$



$$E[\text{mpg, weight}] = (24.5, 2600)$$

The centroid of the cloud

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 31

## Multivariate Expectation

$$E[f(\mathbf{X})] = \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 32



## Test your understanding

Question : When (if ever) does  $E[X + Y] = E[X] + E[Y]$ ?

- All the time?
- Only when X and Y are independent?
- It can fail even if X and Y are independent?

## Bivariate Expectation

$$E[f(x, y)] = \int f(x, y) p(x, y) dy dx$$

$$\text{if } f(x, y) = x \text{ then } E[f(X, Y)] = \int x p(x, y) dy dx$$

$$\text{if } f(x, y) = y \text{ then } E[f(X, Y)] = \int y p(x, y) dy dx$$

$$\text{if } f(x, y) = x + y \text{ then } E[f(X, Y)] = \int (x + y) p(x, y) dy dx$$

$$E[X + Y] = E[X] + E[Y]$$

## Bivariate Covariance

$$\mathbf{s}_{xy} = \text{Cov}[X, Y] = E[(X - \mathbf{m}_x)(Y - \mathbf{m}_y)]$$

$$\mathbf{s}_{xx} = \mathbf{s}_x^2 = \text{Cov}[X, X] = \text{Var}[X] = E[(X - \mathbf{m}_x)^2]$$

$$\mathbf{s}_{yy} = \mathbf{s}_y^2 = \text{Cov}[Y, Y] = \text{Var}[Y] = E[(Y - \mathbf{m}_y)^2]$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 35

## Bivariate Covariance

$$\mathbf{s}_{xy} = \text{Cov}[X, Y] = E[(X - \mathbf{m}_x)(Y - \mathbf{m}_y)]$$

$$\mathbf{s}_{xx} = \mathbf{s}_x^2 = \text{Cov}[X, X] = \text{Var}[X] = E[(X - \mathbf{m}_x)^2]$$

$$\mathbf{s}_{yy} = \mathbf{s}_y^2 = \text{Cov}[Y, Y] = \text{Var}[Y] = E[(Y - \mathbf{m}_y)^2]$$

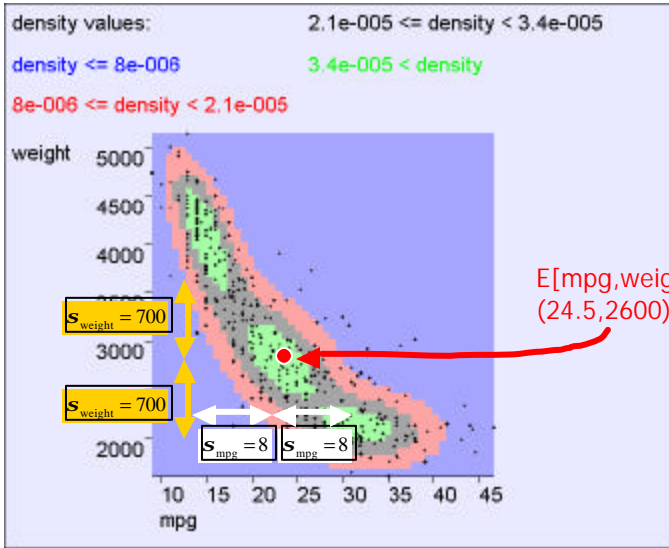
Write  $\mathbf{X} = \begin{pmatrix} X \\ Y \end{pmatrix}$ , then

$$\text{Cov}[\mathbf{X}] = E[(\mathbf{X} - \boldsymbol{\mu}_x)(\mathbf{X} - \boldsymbol{\mu}_x)^T] = \mathbf{S} = \begin{pmatrix} \mathbf{s}_x^2 & \mathbf{s}_{xy} \\ \mathbf{s}_{xy} & \mathbf{s}_y^2 \end{pmatrix}$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 36

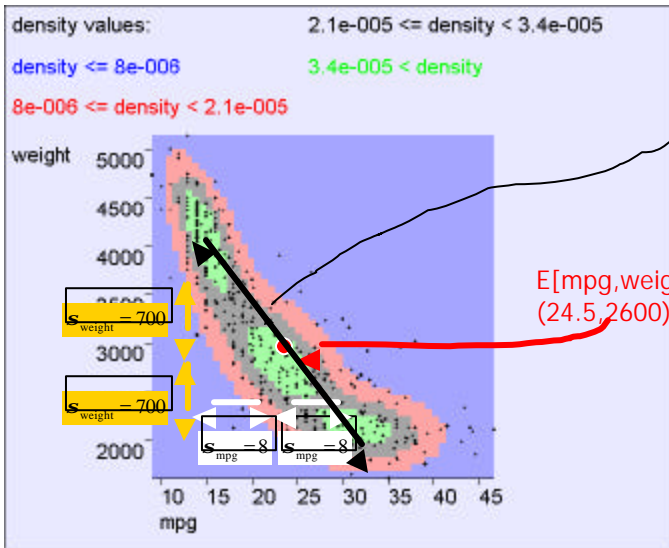
# Covariance Intuition



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 37

# Covariance Intuition



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 38

## Covariance Fun Facts

$$\mathbf{Cov}[\mathbf{X}] = E[(\mathbf{X} - \boldsymbol{\mu}_x)(\mathbf{X} - \boldsymbol{\mu}_x)^T] = \mathbf{S} = \begin{pmatrix} \mathbf{s}_{xx}^2 & \mathbf{s}_{xy} \\ \mathbf{s}_{xy} & \mathbf{s}_{yy}^2 \end{pmatrix}$$

- True or False: If  $\sigma_{xy} = 0$  then X and Y are independent
- True or False: If X and Y are independent then  $\sigma_{xy} = 0$
- True or False: If  $\sigma_{xy} = \sigma_x \sigma_y$  then X and Y are deterministically related
- True or False: If X and Y are deterministically related then  $\sigma_{xy} = \sigma_x \sigma_y$

How could you prove or disprove these?

## General Covariance

Let  $\mathbf{X} = (X_1, X_2, \dots, X_k)$  be a vector of  $k$  continuous random variables

$$\mathbf{Cov}[\mathbf{X}] = E[(\mathbf{X} - \boldsymbol{\mu}_x)(\mathbf{X} - \boldsymbol{\mu}_x)^T] = \mathbf{S}$$

$$\mathbf{S}_{ij} = \text{Cov}[X_i, X_j] = \mathbf{s}_{x_i x_j}$$

S is a  $k \times k$  symmetric non-negative definite matrix

If all distributions are linearly independent it is positive definite

If the distributions are linearly dependent it has determinant zero

## Test your understanding

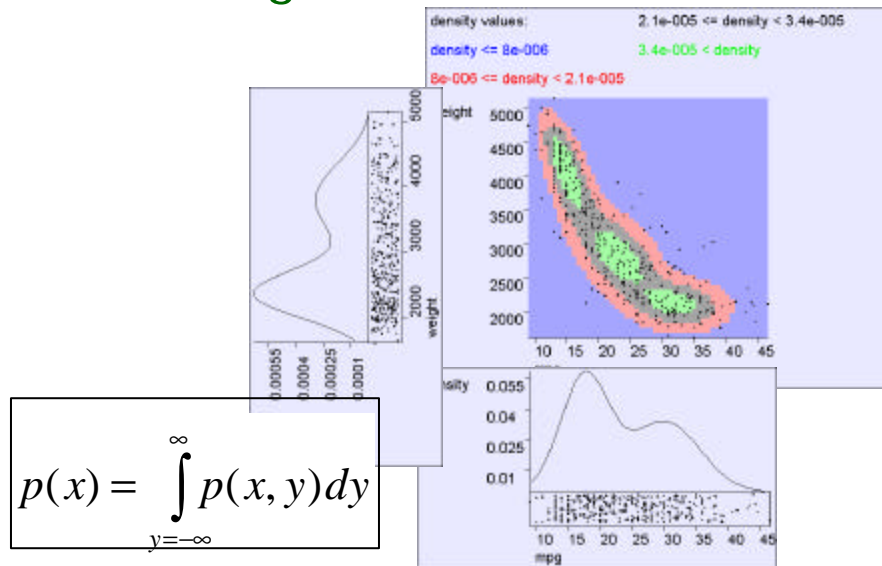
Question : When (if ever) does  $Var[X + Y] = Var[X] + Var[Y]$ ?

- All the time?
- Only when X and Y are independent?
- It can fail even if X and Y are independent?

Copyright © 2001, Andrew W. Moore

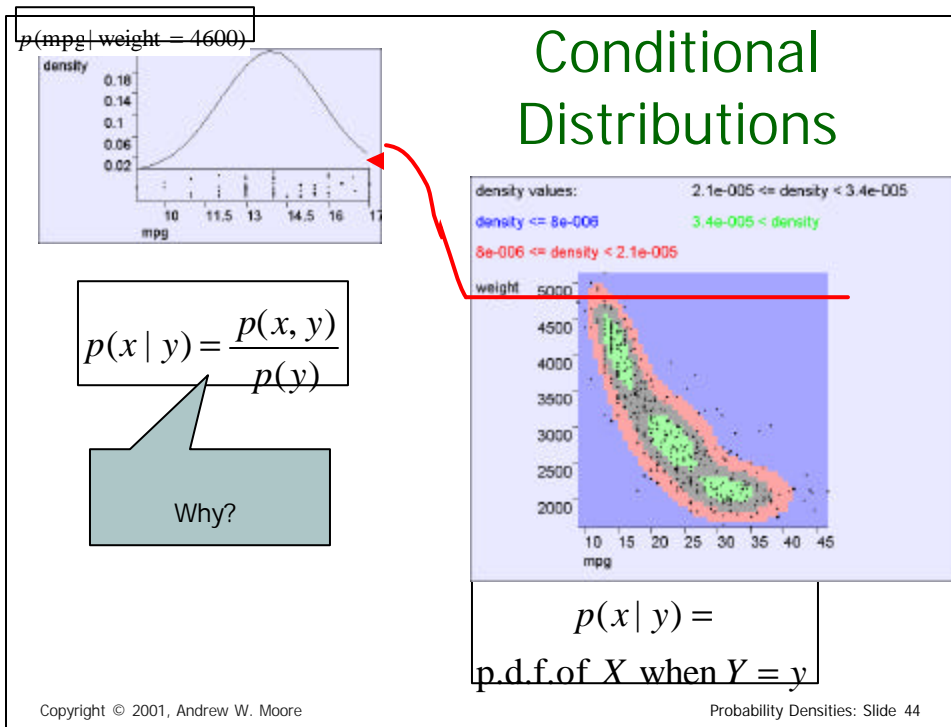
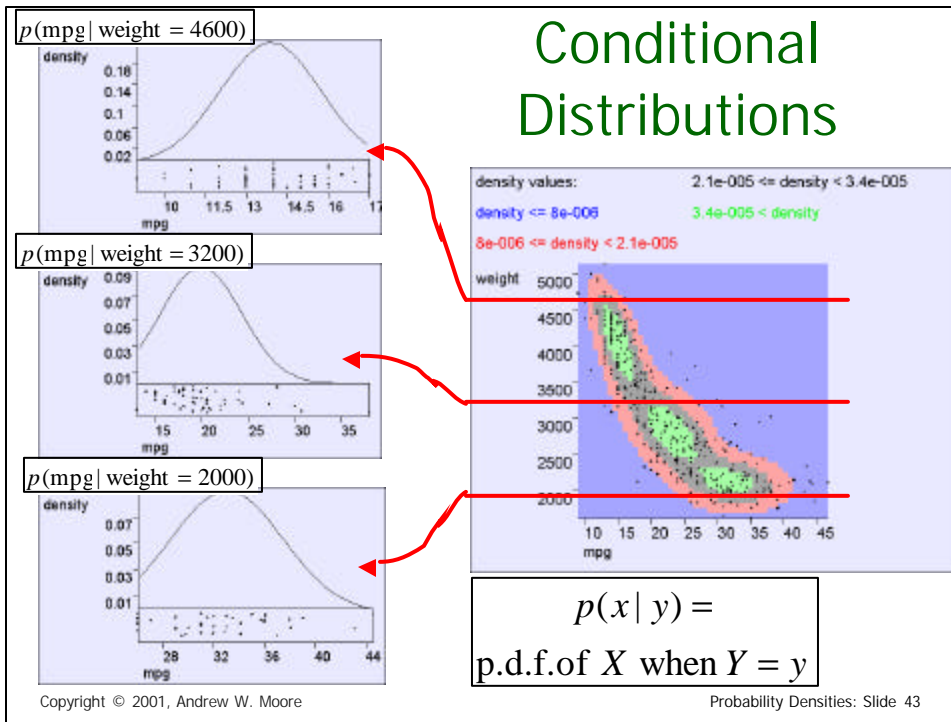
Probability Densities: Slide 41

## Marginal Distributions



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 42



## Independence Revisited

$$X \perp Y \text{ iff } \forall x, y: p(x, y) = p(x)p(y)$$

It's easy to prove that these statements are equivalent...

$$\forall x, y: p(x, y) = p(x)p(y)$$

$$\Leftrightarrow$$

$$\forall x, y: p(x | y) = p(x)$$

$$\Leftrightarrow$$

$$\forall x, y: p(y | x) = p(y)$$

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 45

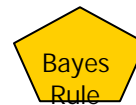
## More useful stuff

$$\int_{x=-\infty}^{\infty} p(x | y) dx = 1$$

(These can all be proved from definitions on previous slides)

$$p(x | y, z) = \frac{p(x, y | z)}{p(y | z)}$$

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



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 46

## Mixing discrete and continuous variables

$$p(x, A = v) = \lim_{h \rightarrow 0} \frac{P\left(x - \frac{h}{2} < X \leq x + \frac{h}{2} \wedge A = v\right)}{h}$$

$$\sum_{v=1}^{n_A} \int_{x=-\infty}^{\infty} p(x, A = v) dx = 1$$

$$p(x | A) = \frac{P(A | x) p(x)}{P(A)}$$

Bayes Rule

$$P(A | x) = \frac{p(x | A) P(A)}{p(x)}$$

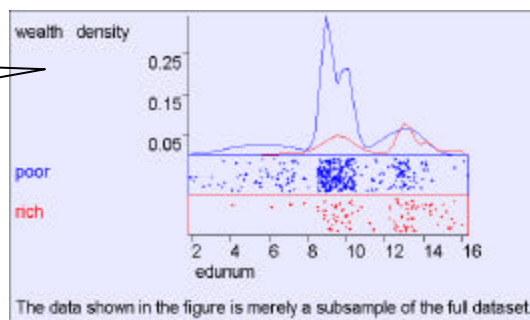
Bayes Rule

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 47

## Mixing discrete and continuous variables

$P(\text{EduYears.Wealthy})$

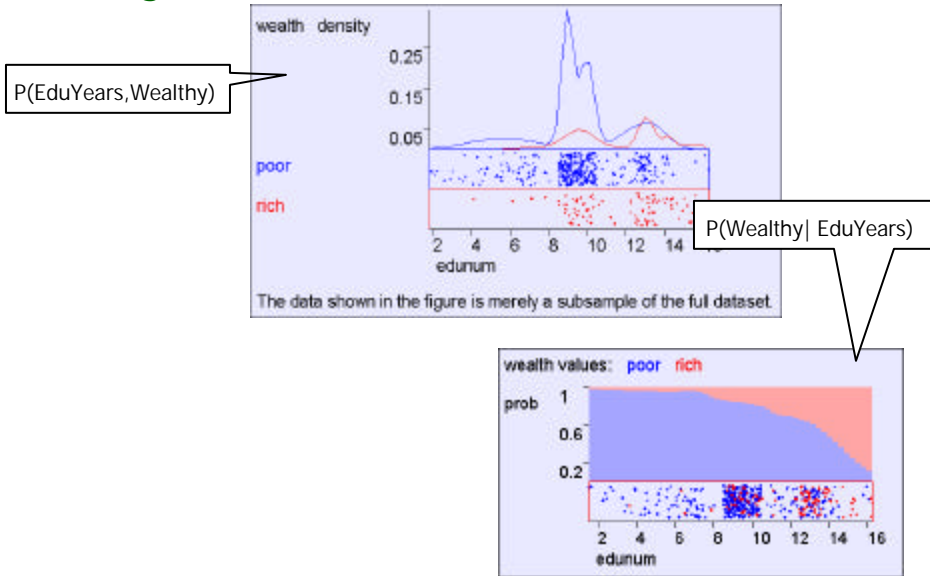


Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 48



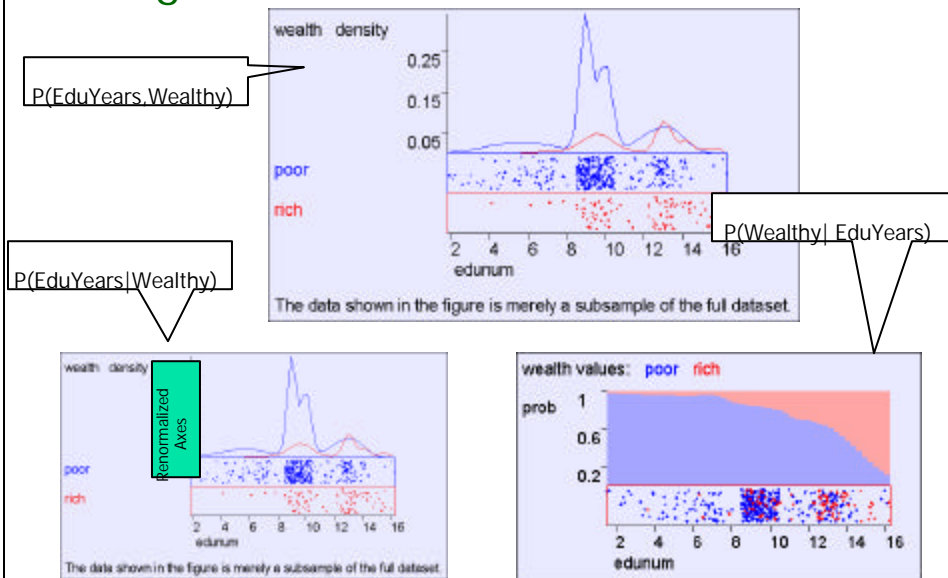
## Mixing discrete and continuous variables



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 49

## Mixing discrete and continuous variables



Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 50

## What you should know

- You should be able to play with discrete, continuous and mixed joint distributions
- You should be happy with the difference between  $p(x)$  and  $P(A)$
- You should be intimate with expectations of continuous and discrete random variables
- You should smile when you meet a covariance matrix
- Independence and its consequences should be second nature

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 51

## Discussion

- Are PDFs the only sensible way to handle analysis of real-valued variables?
- Why is covariance an important concept?
- Suppose  $X$  and  $Y$  are independent real-valued random variables distributed between 0 and 1:
  - What is  $p[\min(X,Y)]$ ?
  - What is  $E[\min(X,Y)]$ ?
- Prove that  $E[X]$  is the value  $u$  that minimizes  $E[(X-u)^2]$
- What is the value  $u$  that minimizes  $E[|X-u|]$ ?

Copyright © 2001, Andrew W. Moore

Probability Densities: Slide 52