# Probability Review

## Lyle Ungar

**I'm comfortable working with probability density functions**

A) yes
B) no

**MLE/MAP**
    Bernoulli and beta distributions
**Probability density functions**
    Gaussians
    Expected value

Yes or No?

Yes

No

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

- **Lectures are recorded**
  - but it's much better to come to class
- **Piazza rocks!**
- **Questions? (chat window)**

**I know the difference between MLE and MAP**
A) yes
B) no

Yes or No?

Yes

No

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

# MLE/MAP

- **MLE maximizes what?**

- **MAP maximizes what?**

- **When is MLE the same as MAP?**

- **We will almost always use MAP. Why?**

# Questions?

**Top**

Start the presentation to see live content. For screen share software, share the entire screen. Get help at **pollev.com/app**

Slide 4

# Probability Densities (PDFs)

## Originally by
## Andrew W. Moore

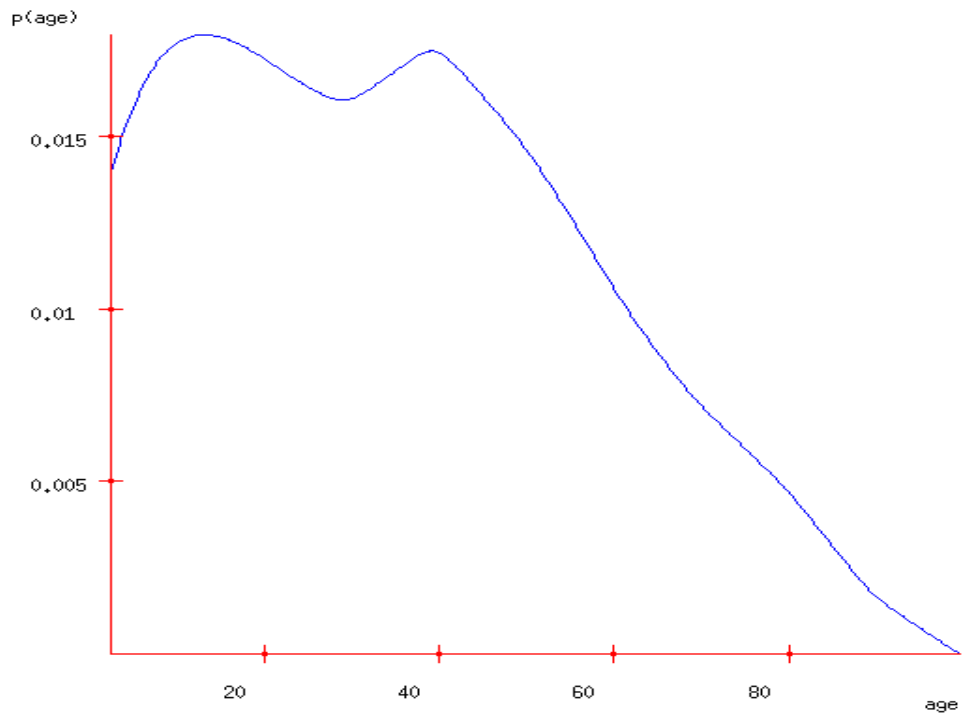Heavily edited by Lyle Ungar

# Probability Densities in ML

- **Why we should care**

- **Notation and fundamentals of continuous PDFs**

- **Multivariate continuous PDFs**

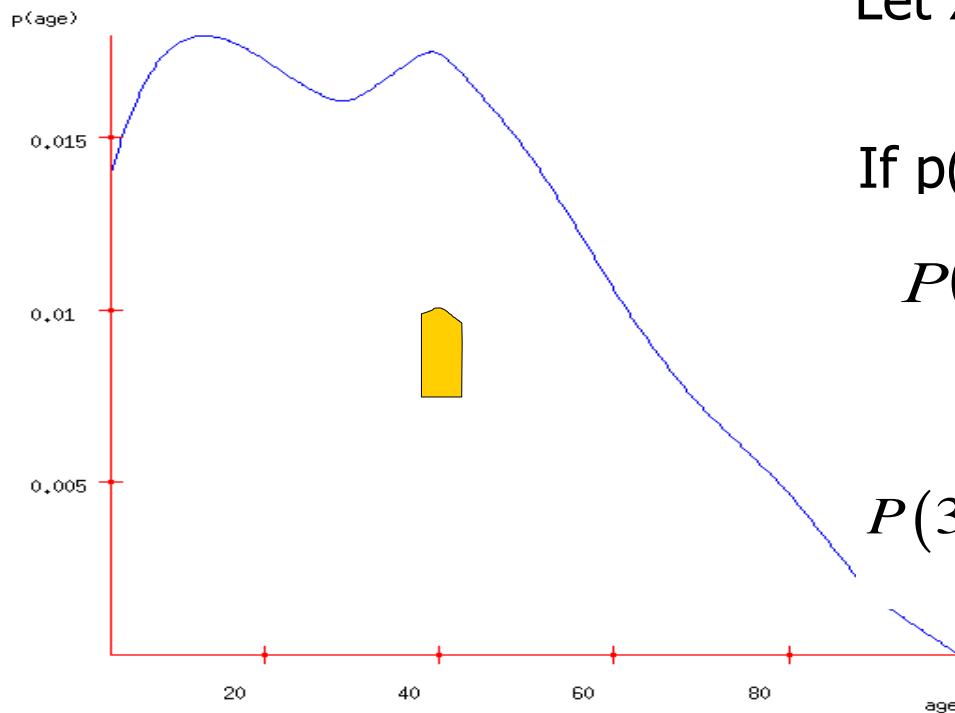- **Expected value, variance, covariance**

# Why we should care

- **Real numbers occur most real data**
  - Can't always quantize them
- **Parameters in models are real valued**
- **You'll need to *intimately* understand PDFs for**
  - kernel methods,
  - clustering with mixture models
  - time series, HMMs
  - proofs about regression

# A PDF of American Ages in 2000

# A PDF of American Ages in 2000

Let X be a continuous random variable.

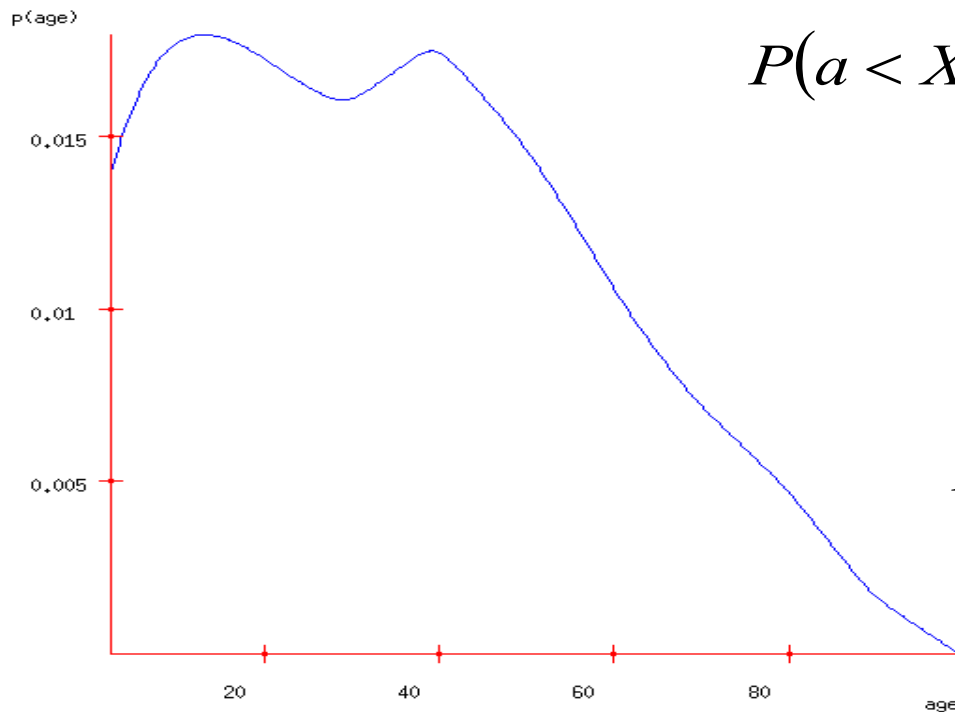If p(x) is a Probability Density

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

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

= 0.36

# Properties of PDFs

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

That means…

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

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

# Properties of PDFs

p(age)

0.015

0.01

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

Therefore...

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

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

80

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 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**

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

**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**

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

**then**

when a value X is sampled from the distribution, you are $\alpha$ 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**

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

**then**

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

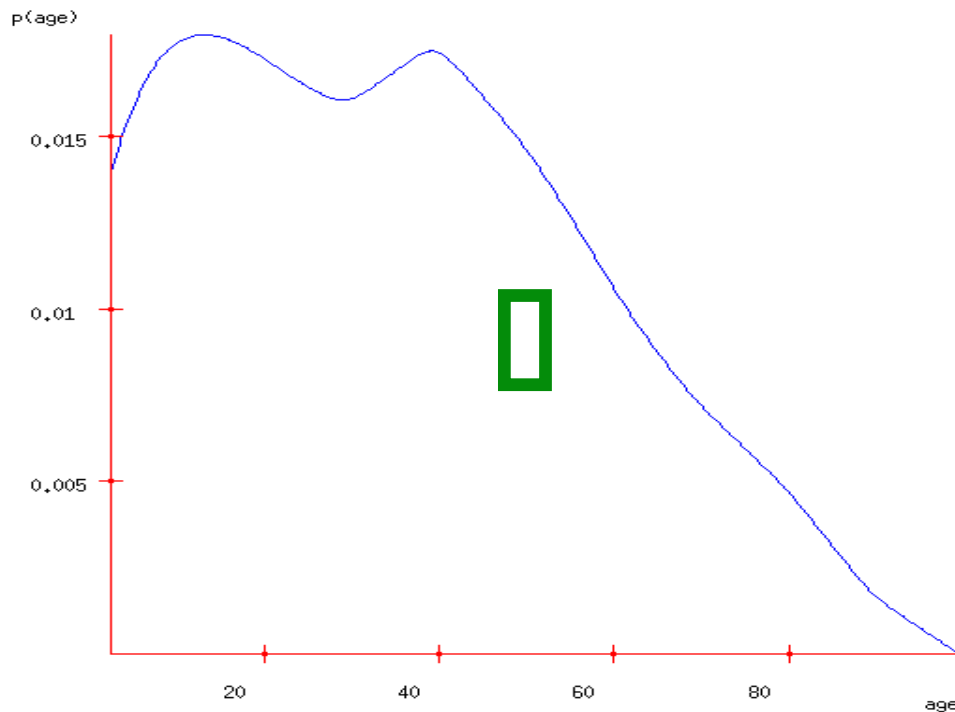# Talking to your stomach

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

**If**

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

**then**

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

# 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(age) stop and return age

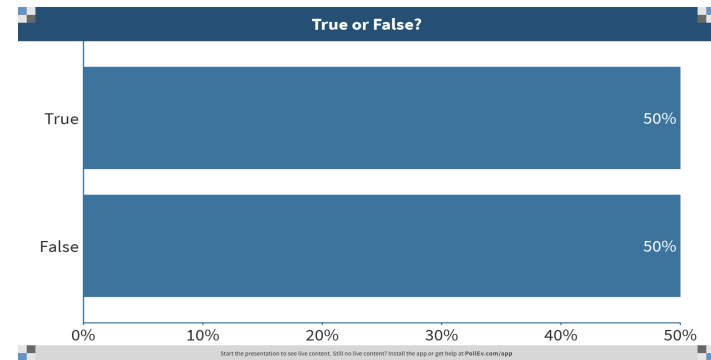3.  Else try again: go to Step 1.

# Test your understanding
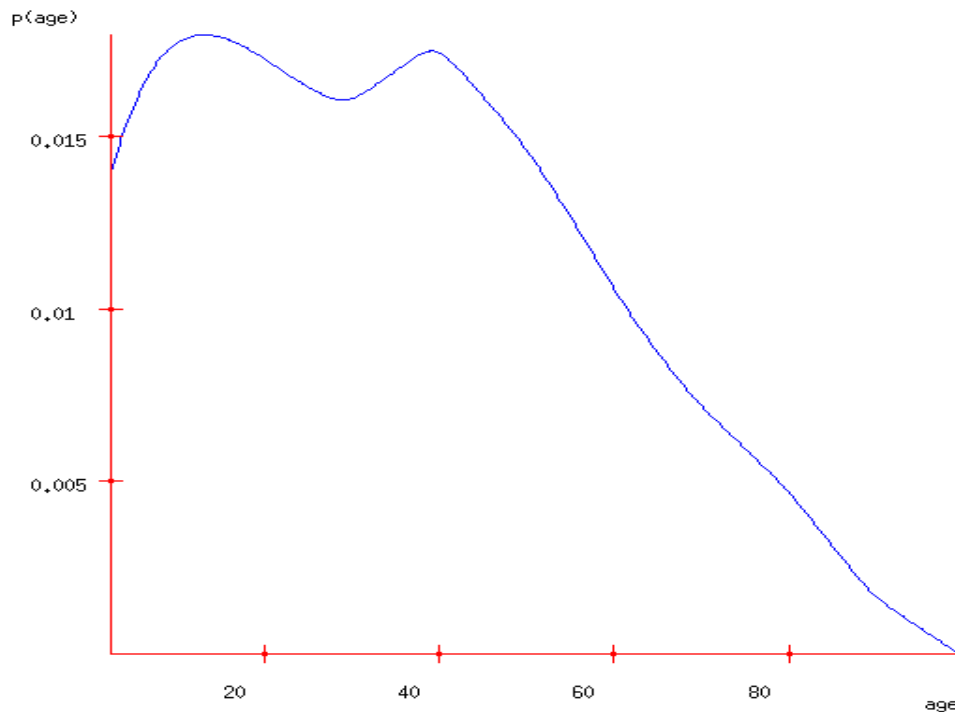
- **True or False:**

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

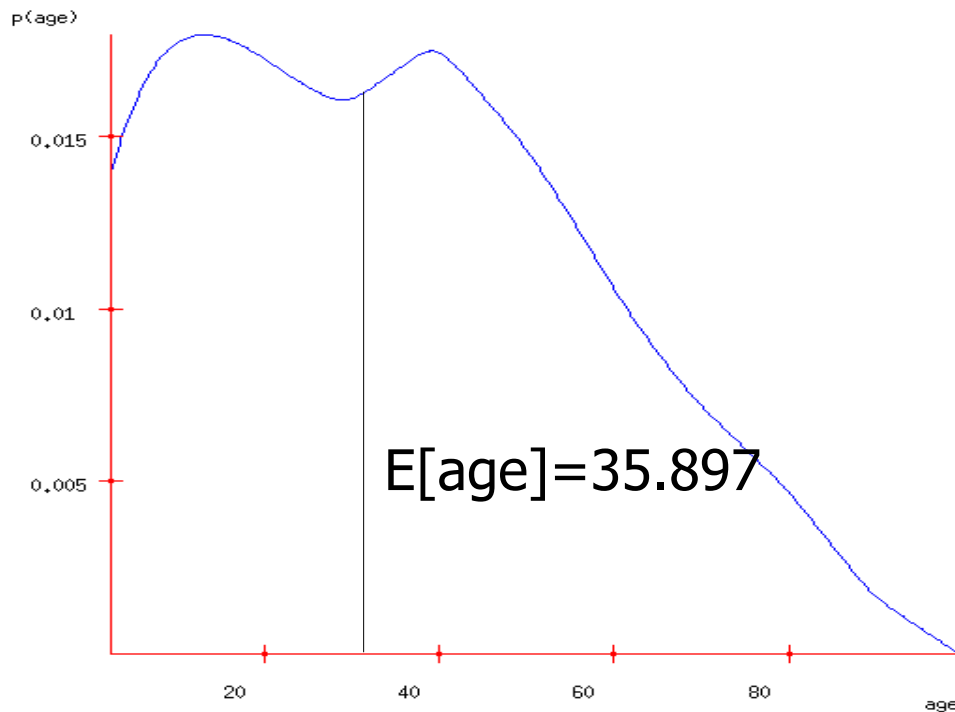- **True or False:**

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

# 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

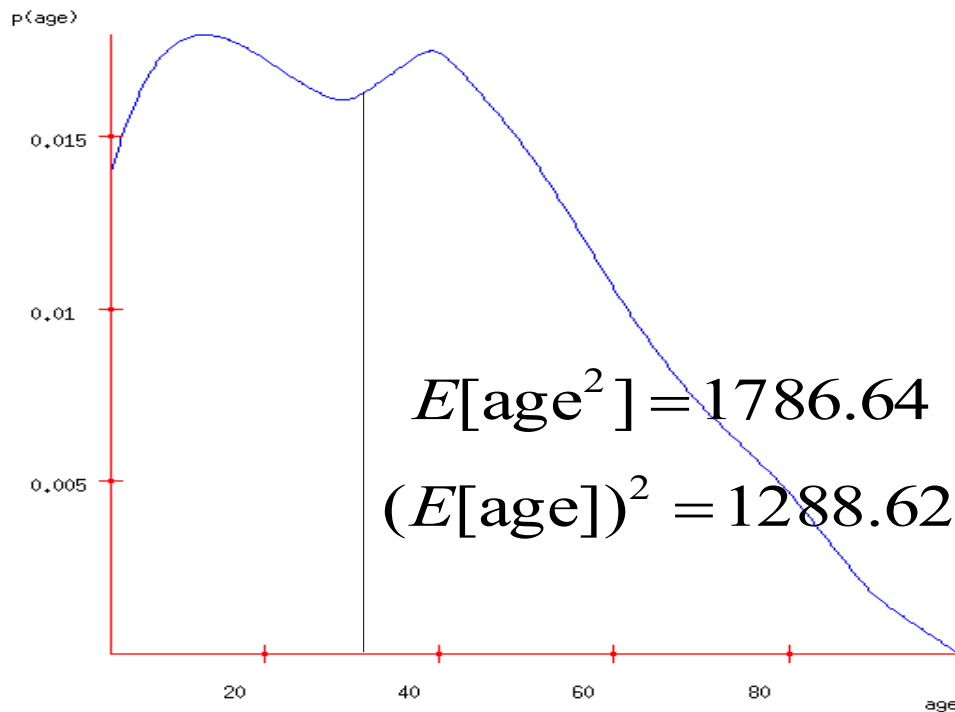$$= \int_{x=-\infty}^{\infty} x \, p(x) \, dx$$

# Expectations



E[age]=35.897

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

= 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

# Expectation of a function



$$E[\text{age}^2] = 1786.64$$
$$(E[\text{age}])^2 = 1288.62$$

$\mu$=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)

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

Note that in general:

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

# Variance

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

$$\sigma^2 = \int\limits_{x=-\infty}^{\infty} (x-\mu)^2 \, p(x)\,dx$$

$$\mathrm{Var[age]} = 498.02$$

= 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

# Standard Deviation
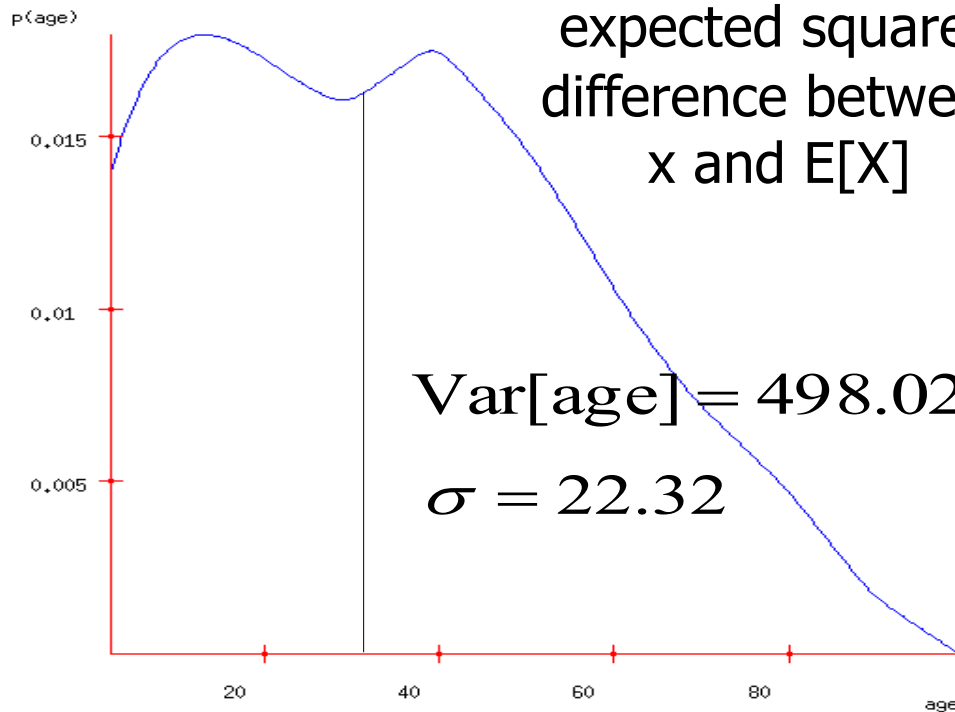
$\sigma^2$ = Var[X] = the expected squared difference between x and E[X]

$$\sigma^2 = \int\limits_{x=-\infty}^{\infty} (x-\mu)^2 \, p(x) \, dx$$



$$\mathrm{Var[age]} = 498.02$$

$$\sigma = 22.32$$

= 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
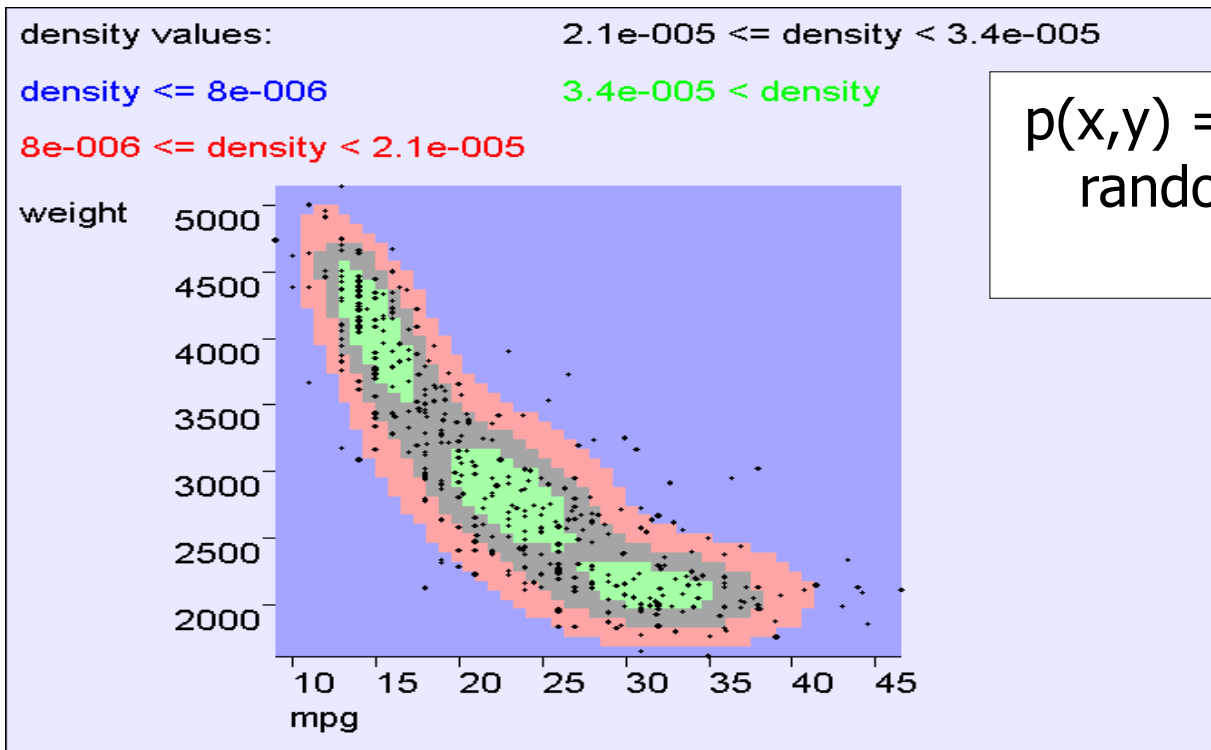
$\sigma$ = Standard Deviation = "typical" deviation of X from its mean $\sigma = \sqrt{\mathrm{Var[X]}}$

# In two dimensions

density values:          2.1e-005 <= density < 3.4e-005

density <= 8e-006        3.4e-005 < density

8e-006 <= density < 2.1e-005
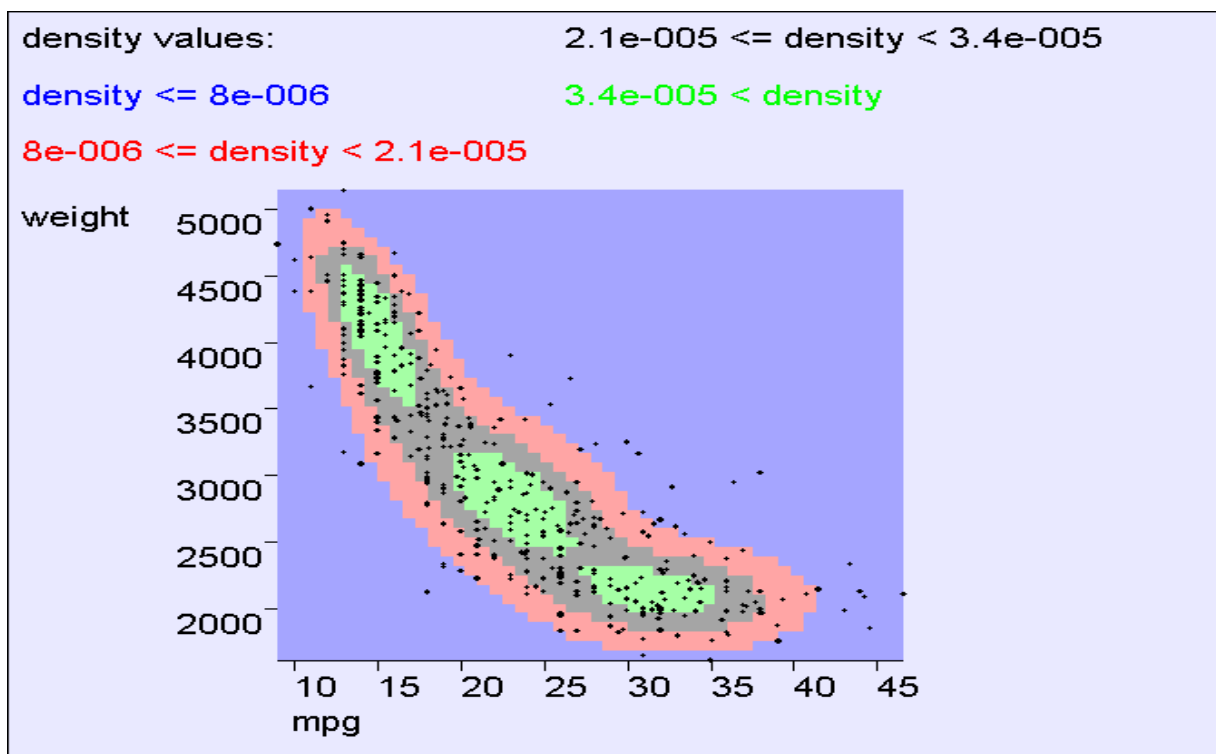


p(x,y) = probability density of random variables (X,Y) at location (x,y)

# In two dimensions

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

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

density values:      2.1e-005 <= density < 3.4e-005

density <= 8e-006      3.4e-005 < density

8e-006 <= density < 2.1e-005

# In two dimensions

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

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

density values:      2.1e-005 <= density < 3.4e-005

density <= 8e-006      3.4e-005 < density

8e-006 <= density < 2.1e-005

P( 20<mpg<30 and 2500<weight<3000) =

area under the 2-d surface within the red rectangle

weight

5000
4500
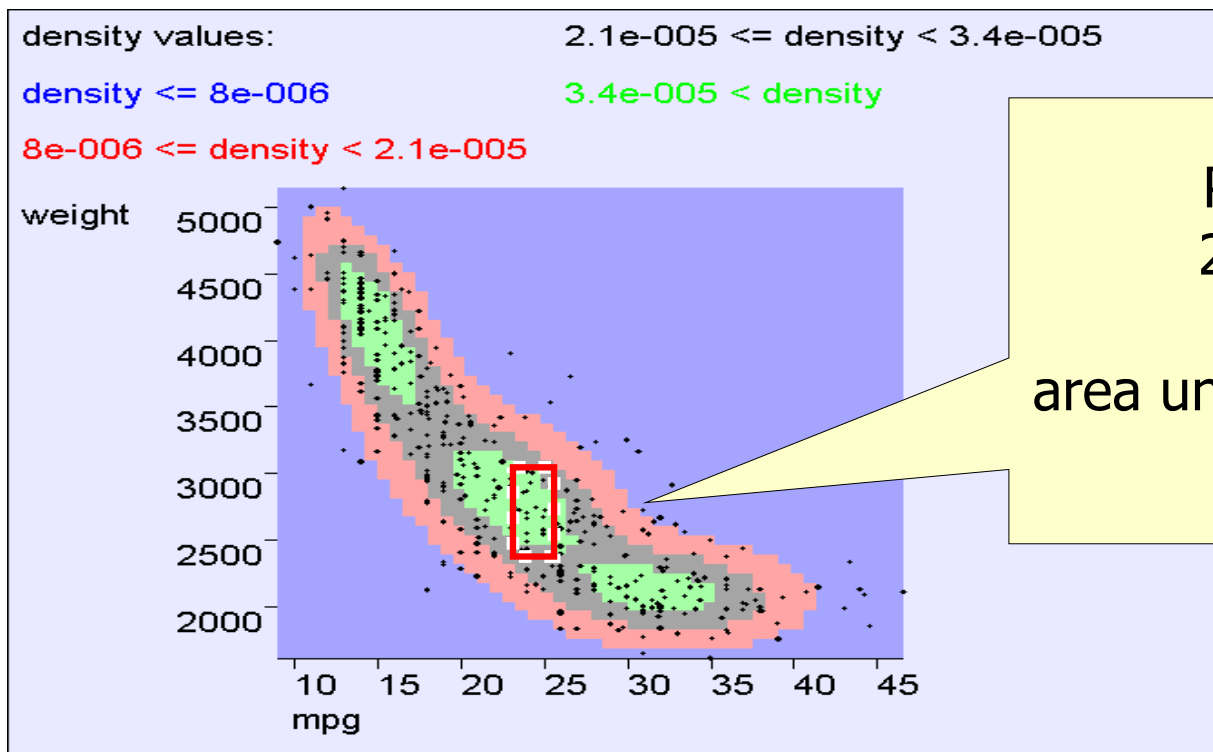4000
3500
3000
2500
2000

10  15  20  25  30  35  40  45
mpg

# In two dimensions

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

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

density values:

density <= 8e-006

8e-006 <= density < 2.1e-005

2.1e-005 <= density < 3.4e-005

3.4e-005 < density

weight

5000
4500
4000
3500
3000
2500
2000

10  15  20  25  30  35  40  45
mpg

P( [(mpg-25)/10]² + [(weight-3300)/1500]² < 1 ) =

area under the 2-d surface within the red oval

# In two dimensions

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

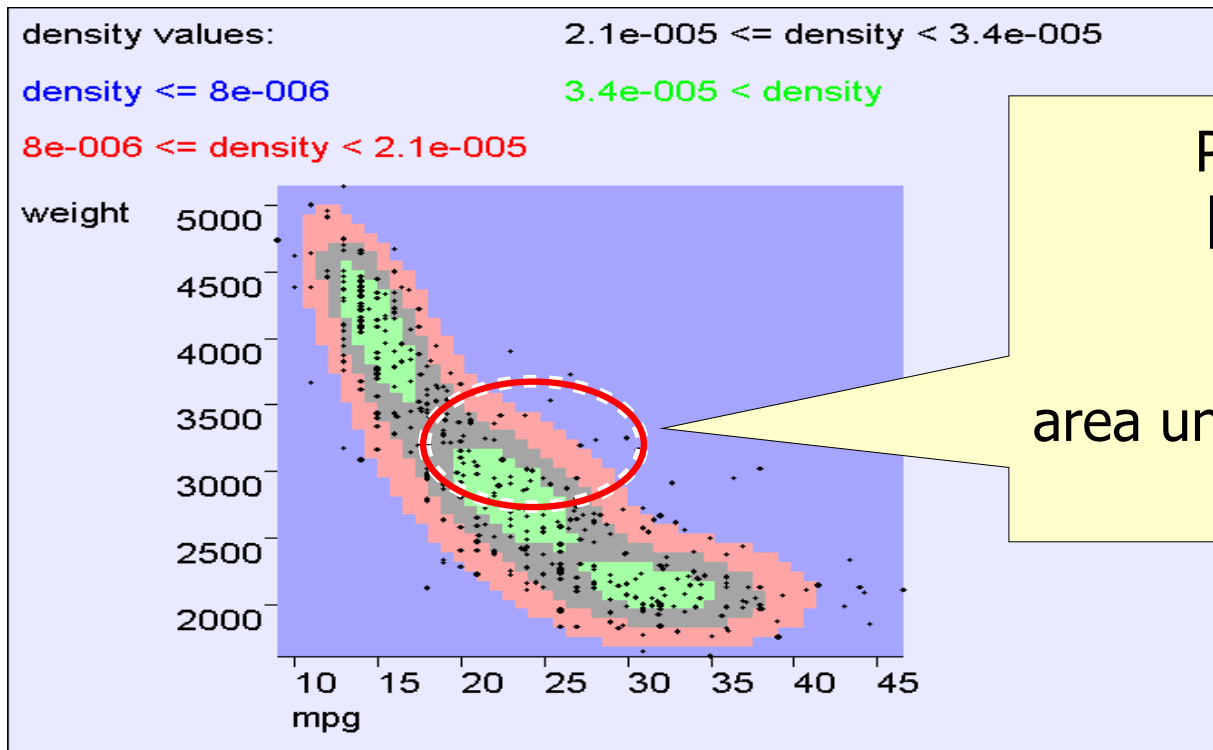$$P((X,Y) \in R) = \iint\limits_{(x,y) \in R} p(x,y) \, dy \, dx$$

Take the special case of region R = "everywhere".

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

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

# In two dimensions

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

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

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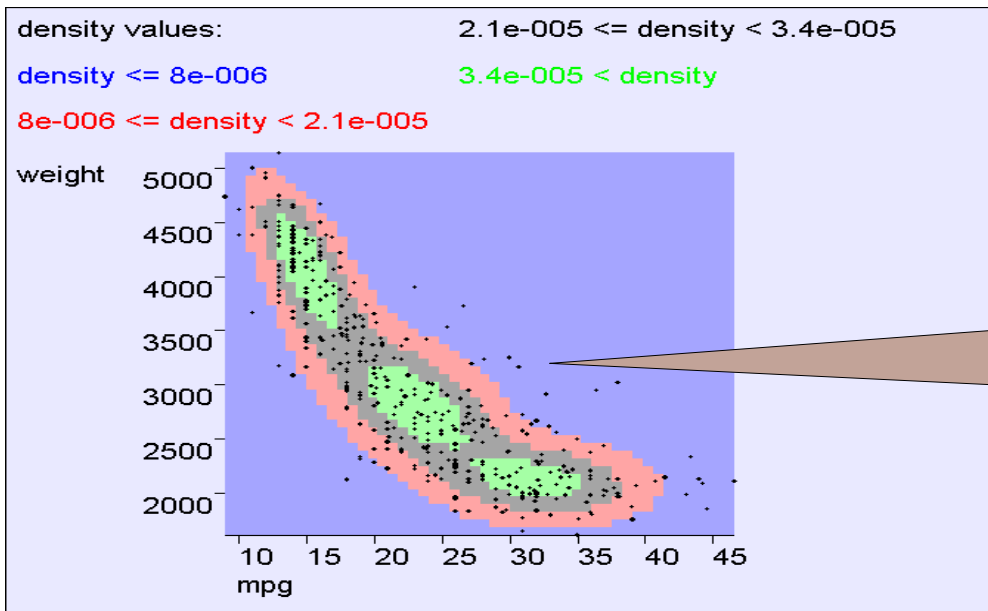
## In m dimensions

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

$$P((X_1, X_2, ..., X_m) \in R) =$$

$$\iint ... \int_{(x_1, x_2, ..., x_m) \in R} p(x_1, x_2, ..., x_m) dx_m, ... dx_2, dx_1$$

# Independence

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

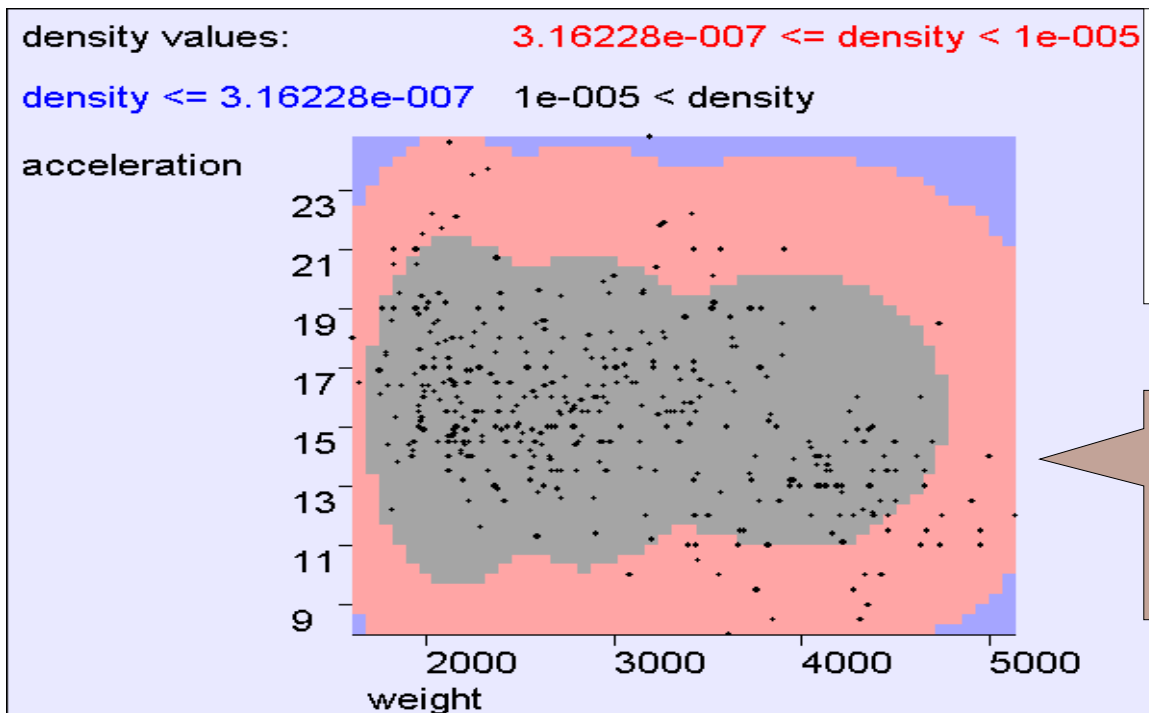If X and Y are independent then knowing the value of X does not help predict the value of Y

density values:        2.1e-005 <= density < 3.4e-005

density <= 8e-006        3.4e-005 < density

8e-006 <= density < 2.1e-005



mpg,weight NOT independent

# Independence

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



density values:  3.16228e-007 <= density < 1e-005

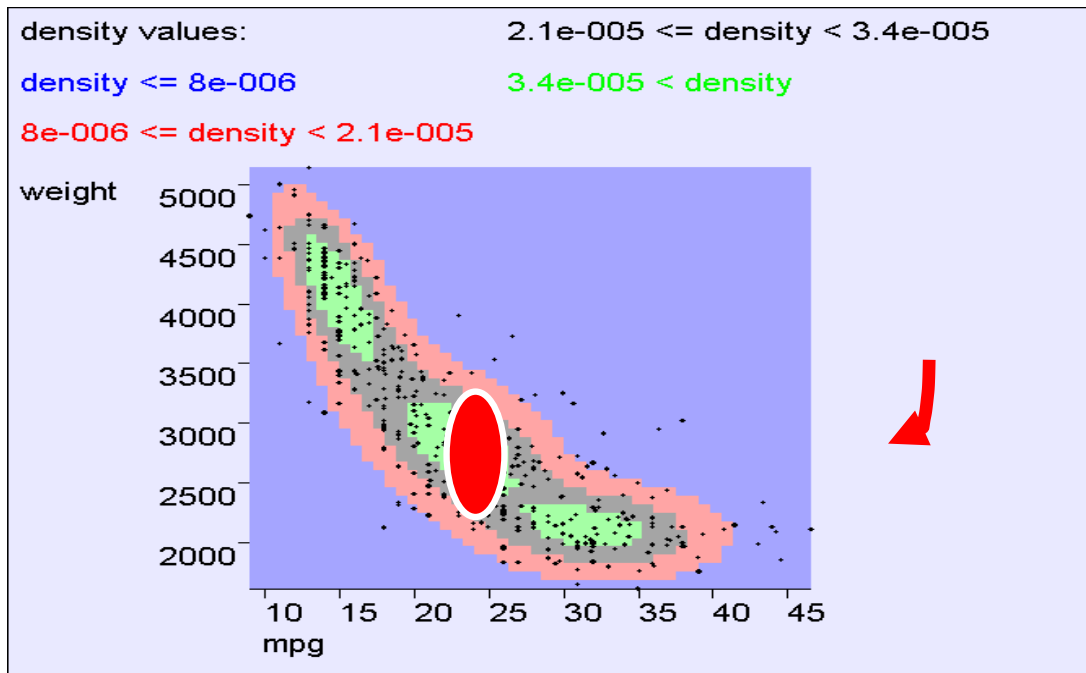density <= 3.16228e-007    1e-005 < density

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

# Multivariate Expectation

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



E[mpg,weight] =
(24.5,2600)

The centroid of the cloud

# Multivariate Expectation

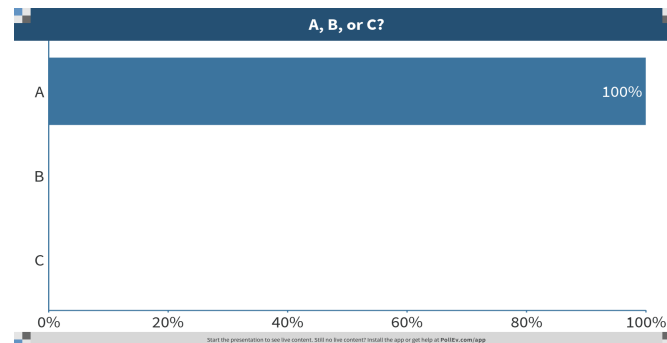$$E[f(\mathbf{X})] = \int f(\mathbf{x}) \ p(\mathbf{x}) d\mathbf{x}$$

# Test your understanding

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

A) All the time?

B) Only when X and Y are independent?

C) It can fail even if X and Y are independent?



A, B, or C?

| | |
|---|---|
| A | 100% |
| B | |
| C | |

0%    20%    40%    60%    80%    100%

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

# Bivariate Expectation

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

if $f(x,y) = x$ then $E[f(X,Y)] = \int x \, p(x,y) dydx$

if $f(x,y) = y$ then $E[f(X,Y)] = \int y \, p(x,y) dydx$

if $f(x,y) = x+y$ then $E[f(X,Y)] = \int (x+y) \, p(x,y) dydx$

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

# Bivariate Covariance

$$\sigma_{xy} = \text{Cov}[X,Y] = E[(X - \mu_x)(Y - \mu_y)]$$

$$\sigma_{xx} = \sigma^2{}_x = \text{Cov}[X,X] = Var[X] = E[(X - \mu_x)^2]$$

$$\sigma_{yy} = \sigma^2{}_y = \text{Cov}[Y,Y] = Var[Y] = E[(Y - \mu_y)^2]$$

# Bivariate Covariance

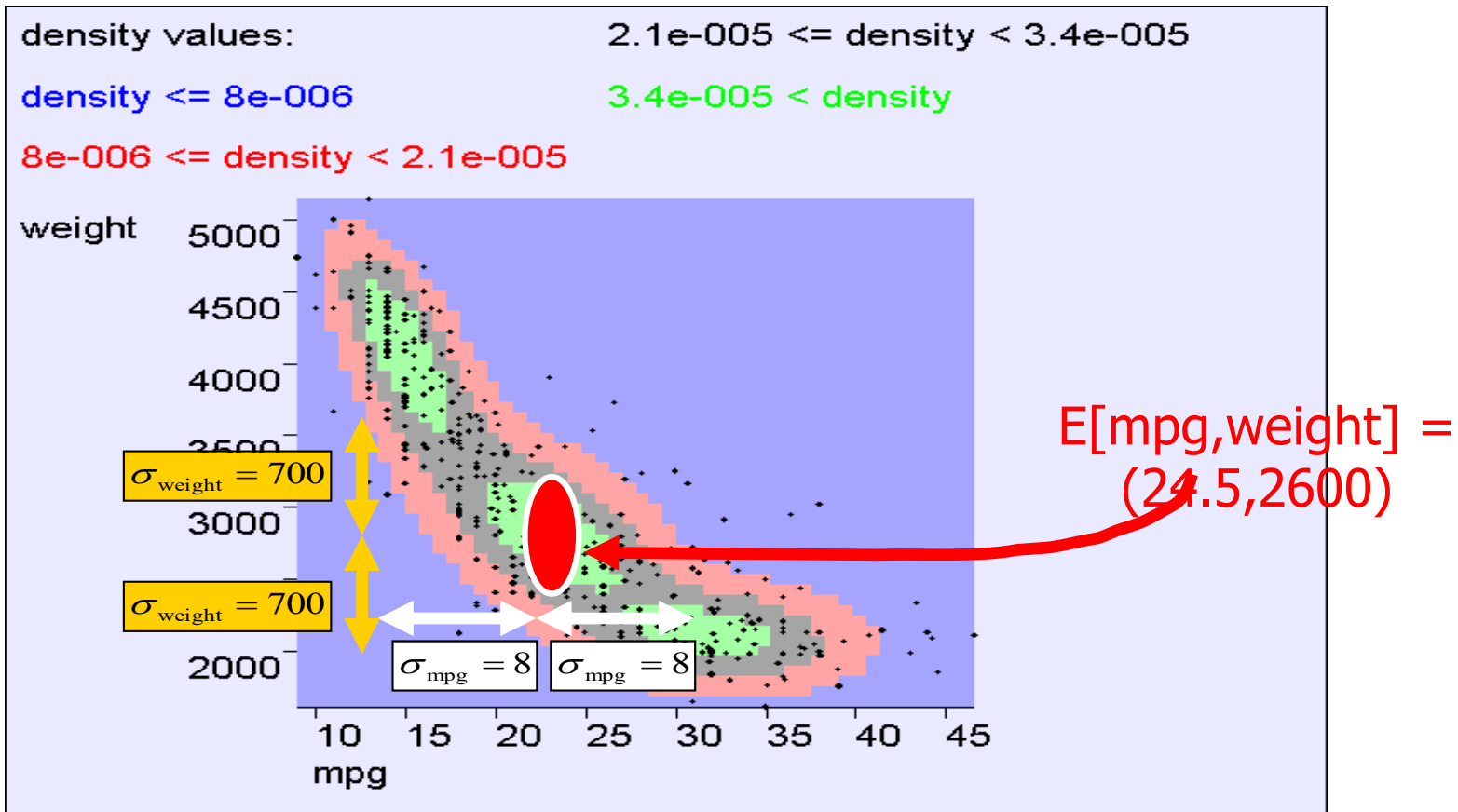$$\sigma_{xy} = \text{Cov}[X, Y] = E[(X - \mu_x)(Y - \mu_y)]$$

$$\sigma_{xx} = \sigma^2_x = \text{Cov}[X, X] = Var[X] = E[(X - \mu_x)^2]$$

$$\sigma_{yy} = \sigma^2_y = \text{Cov}[Y, Y] = Var[Y] = E[(Y - \mu_y)^2]$$

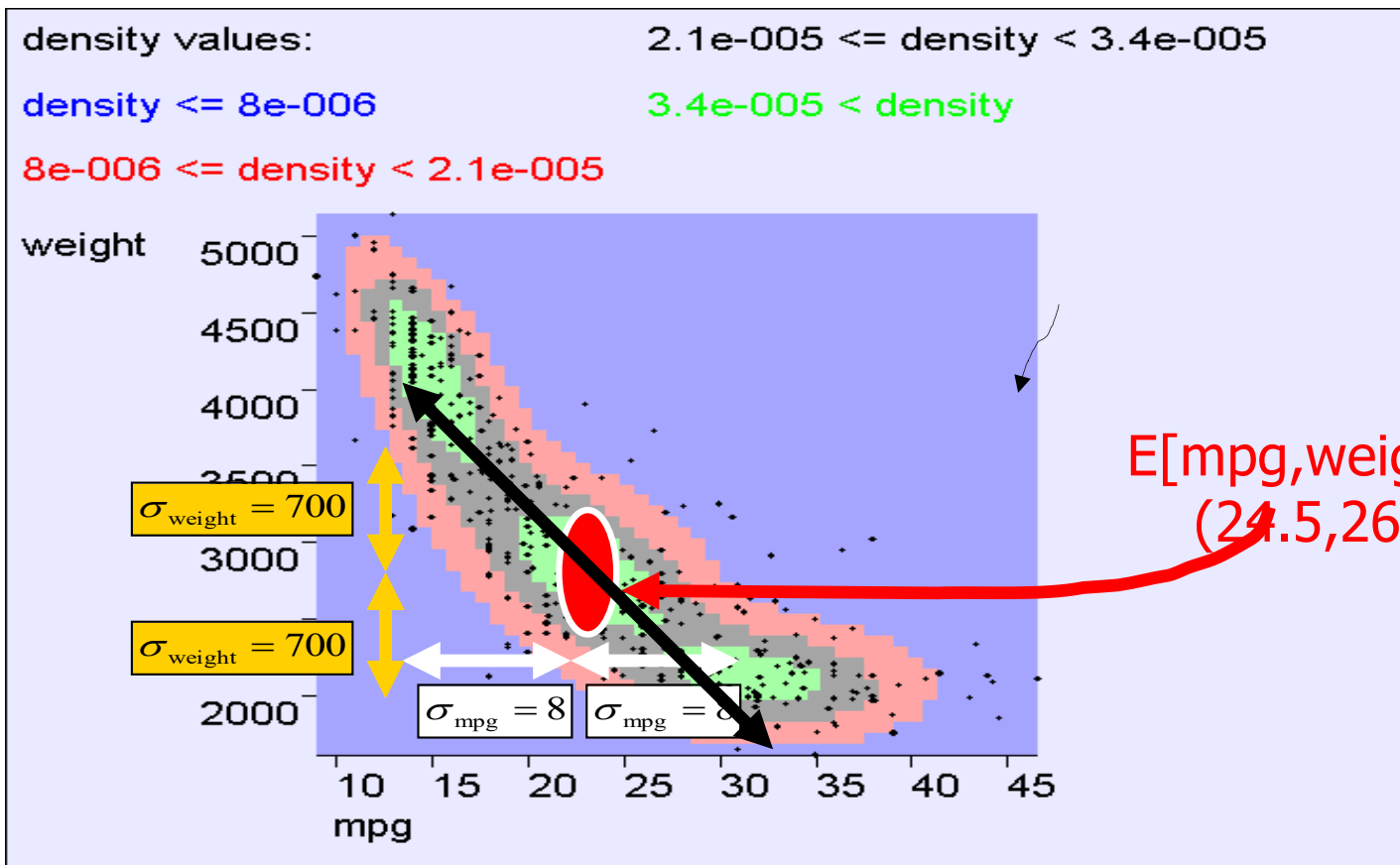$$\text{Write } \mathbf{X} = \begin{pmatrix} X \\ Y \end{pmatrix}, \text{ then}$$

$$\mathbf{Cov[\ X]} = E[(\mathbf{X} - \boldsymbol{\mu}_x)(\mathbf{X} - \boldsymbol{\mu}_x)^T] = \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2_x & \sigma_{xy} \\ \sigma_{xy} & \sigma^2_y \end{pmatrix}$$
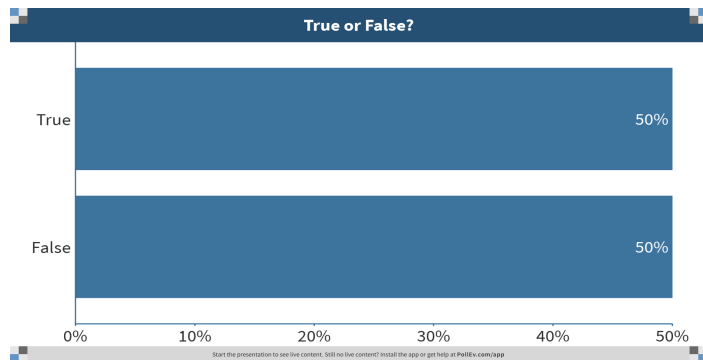
# Covariance Intuition



density values:

density <= 8e-006

8e-006 <= density < 2.1e-005

2.1e-005 <= density < 3.4e-005

3.4e-005 < density

weight

$\sigma_{\text{weight}} = 700$

$\sigma_{\text{weight}} = 700$

$\sigma_{\text{mpg}} = 8$

$\sigma_{\text{mpg}} = 8$

E[mpg,weight] = (24.5,2600)

mpg

# Covariance Intuition



density values:

density <= 8e-006

8e-006 <= density < 2.1e-005

2.1e-005 <= density < 3.4e-005

3.4e-005 < density

weight

5000
4500
4000
3500
3000
2000

$\sigma_{weight} = 700$

$\sigma_{weight} = 700$

$\sigma_{mpg} = 8$

$\sigma_{mpg} = 8$

10 15 20 25 30 35 40 45

mpg

Principal
Eigenvector
of $\Sigma$

E[mpg,weight] =
(24.5,2600)

# Covariance Fun Facts

$$\mathbf{Cov[\ X]\ } = E[(\mathbf{X} - \boldsymbol{\mu}_x)(\mathbf{X} - \boldsymbol{\mu}_x)^T] = \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2{}_x & \sigma_{xy} \\ \sigma_{xy} & \sigma^2{}_y \end{pmatrix}$$

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

**True or False?**

| | |
|---|---|
| True | 50% |
| False | 50% |

0%   10%   20%   30%   40%   50%

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

How could you prove or disprove these?

# Covariance Fun Facts

For example, let $X$ be uniformly distributed in $[-1, 1]$ and let $Y = X^2$.

Clearly, $X$ and $Y$ are dependent, but

$$
\begin{aligned}
\mathrm{cov}(X, Y) &= \mathrm{cov}(X, X^2) \\
&= \mathrm{E}[X \cdot X^2] - \mathrm{E}[X] \cdot \mathrm{E}[X^2] \\
&= \mathrm{E}[X^3] - \mathrm{E}[X]\,\mathrm{E}[X^2] \\
&= 0 - 0 \cdot \mathrm{E}[X^2] \\
&= 0.
\end{aligned}
$$

# Covariance Fun Facts

$$\mathbf{Cov}[\,\mathbf{X}\,] = E[(\mathbf{X}-\boldsymbol{\mu}_x)(\mathbf{X}-\boldsymbol{\mu}_x)^T] = \boldsymbol{\Sigma} = \begin{pmatrix} \sigma^2{}_x & \sigma_{xy} \\ \sigma_{xy} & \sigma^2{}_y \end{pmatrix}$$

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

**True or False?**

True      50%

False      50%

0%   10%   20%   30%   40%   50%

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

https://en.wikipedia.org/wiki/Covariance
#Uncorrelatedness_and_independence

# Covariance Fun Facts

$$\mathbf{Cov[\,X\,]} \;=\; E[(\mathbf{X}-\boldsymbol{\mu}_x)(\mathbf{X}-\boldsymbol{\mu}_x)^T\,] \;=\; \boldsymbol{\Sigma} \;=\; \begin{pmatrix} \sigma^2{}_x & \sigma_{xy} \\ \sigma_{xy} & \sigma^2{}_y \end{pmatrix}$$

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

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

$$\boldsymbol{\Sigma}_{ij} = Cov[X_i, X_j] = \sigma_{x_i x_j}$$

S is a k x k symmetric positive semi-definite (PSD) matrix

If all distributions are linearly independent it is positive definite

If the distributions are linearly dependent it has at least on zero eigenvalue

# Test your understanding

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

A) All the time?

B) Only when X and Y are independent?
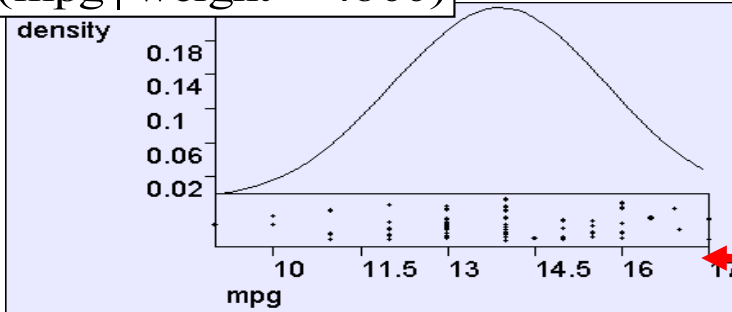
C) It can fail even if X and Y are independent?



A, B, or C?

A — 100%
B
C

0%    20%    40%    60%    80%    100%

Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

# Marginal Distributions



density values:               2.1e-005 <= density < 3.4e-005

density <= 8e-006          3.4e-005 < density

8e-006 <= density < 2.1e-005

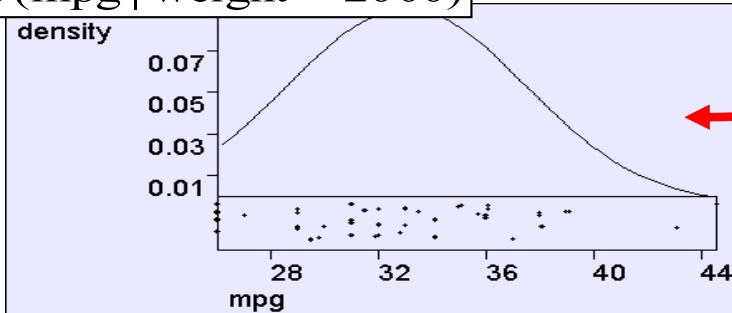$$p(x) = \int\limits_{y=-\infty}^{\infty} p(x, y)\, dy$$

# Conditional Distributions

$p(\text{mpg} \mid \text{weight} = 4600)$

$p(\text{mpg} \mid \text{weight} = 3200)$
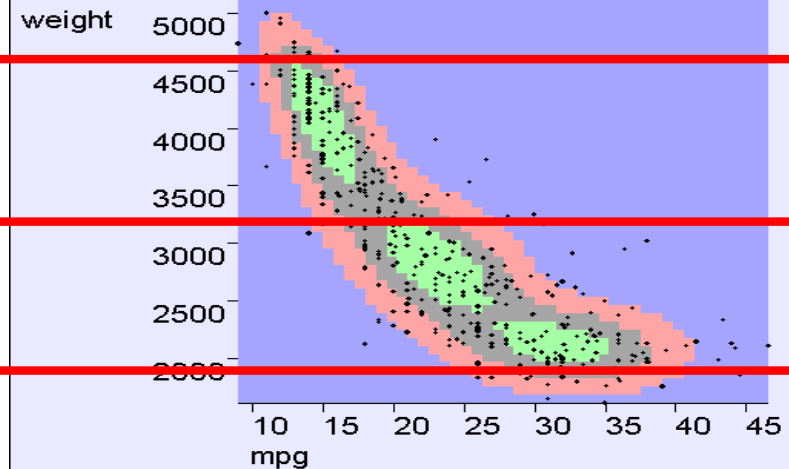
$p(\text{mpg} \mid \text{weight} = 2000)$

density values:

density <= 8e-006
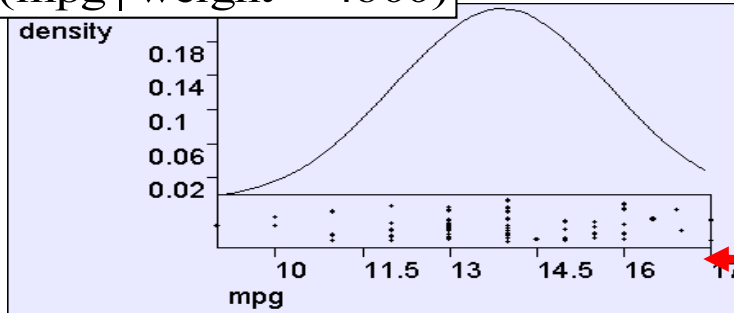
8e-006 <= density < 2.1e-005

2.1e-005 <= density < 3.4e-005

3.4e-005 < density

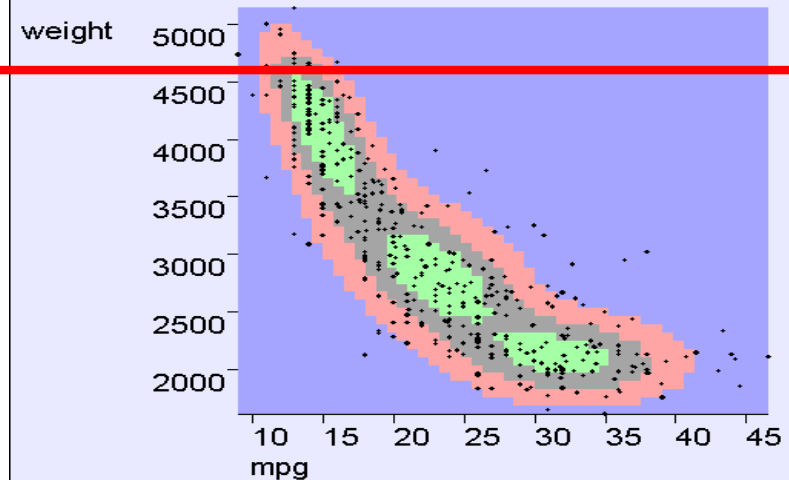$$p(x \mid y) =$$
$$\text{p.d.f. of } X \text{ when } Y = y$$

# Conditional Distributions

$p(\text{mpg} \mid \text{weight} = 4600)$



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

Why?



density values:

density <= 8e-006

8e-006 <= density < 2.1e-005

2.1e-005 <= density < 3.4e-005

3.4e-005 < density

$$p(x \mid y) =$$
p.d.f. of $X$ when $Y = y$

# Independence Revisited

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

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

$$\forall \mathrm{x}, \mathrm{y} : p(x, y) = p(x)p(y)$$

$$\Longleftrightarrow$$

$$\forall \mathrm{x}, \mathrm{y} : p(x \mid y) = p(x)$$

$$\Longleftrightarrow$$

$$\forall \mathrm{x}, \mathrm{y} : p(y \mid x) = p(y)$$

# More useful stuff

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

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

$$p(x \mid y, z) = \frac{p(x, y \mid z)}{p(y \mid z)}$$
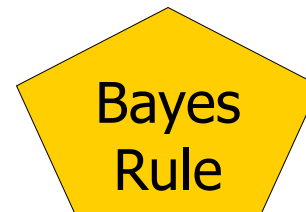
$$p(x \mid y) = \frac{p(y \mid x)\,p(x)}{p(y)}$$
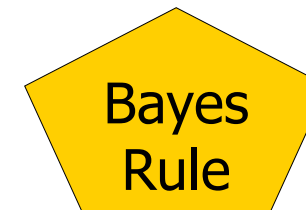
Bayes Rule

# Mixing discrete and continuous variables

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

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

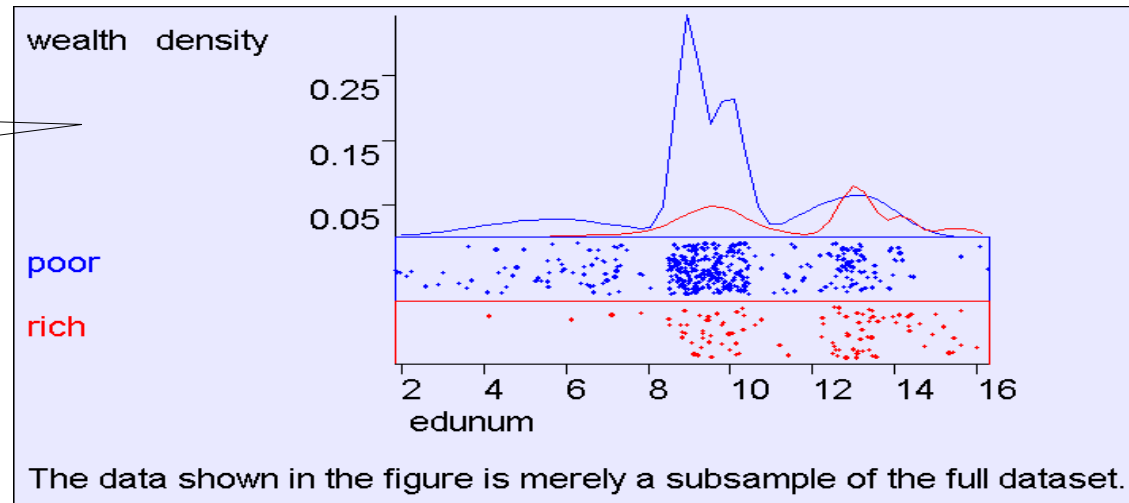$$p(x \mid A) = \frac{P(A \mid x)\, p(x)}{P(A)}$$

Bayes Rule

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

Bayes Rule

# Mixing discrete and continuous variables

P(EduYears,Wealthy)



The data shown in the figure is merely a subsample of the full dataset.

# Mixing discrete and continuous variables



P(EduYears,Wealthy)

P(Wealthy| EduYears)

# Mixing discrete and continuous variables

P(EduYears,Wealthy)



wealth density

0.25
0.15
0.05

poor

rich

2   4   6   8   10   12   14
edunum

The data shown in the figure is merely a subsample of the full dataset.

P(Wealthy| EduYears)

P(EduYears|Wealthy)



wealth density

Renormalized Axes

poor

rich

2   4   6   8   10   12   14   16
edunum

The data shown in the figure is merely a subsample of the full dataset.



wealth values:   poor   rich

prob

1
0.6
0.2

2   4   6   8   10   12   14   16
edunum

# What you should know

- **You should**

  - be able to play with discrete, continuous and mixed joint distributions

  - be happy with the difference between $p(x)$ and $P(A)$

  - be intimate with expectations, variance and covariance of continuous and discrete random variables

  - smile when you meet a covariance matrix

- **Independence and its consequences should be second nature**

# What questions do you have on today's class?

**Top**

Start the presentation to see live content. For screen share software, share the entire screen. Get help at **pollev.com/app**

Slide 58

# How is my speed?

Slow

Good

Fast