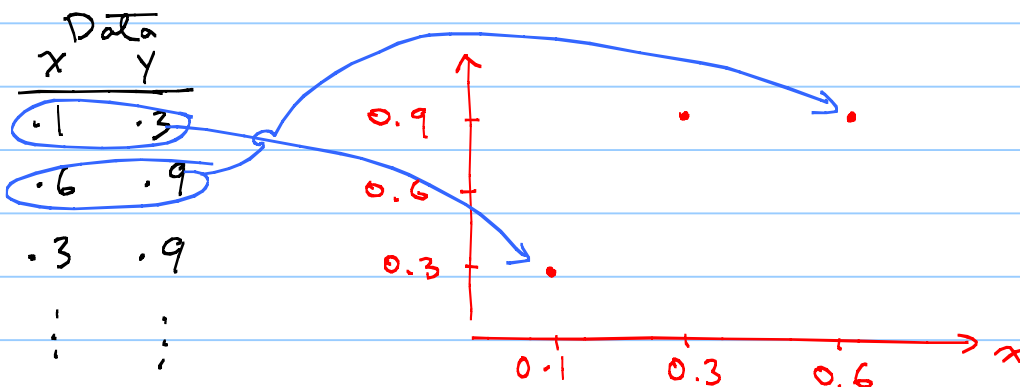# Lecture 12 (Ch.3)

Thus far, our focus has been on 1 column of data, and 1 variable. I.e. univariate analysis.

With 2 (or more) variables, we can do all of the above, but we can also ask about the <u>relationship between</u> them.
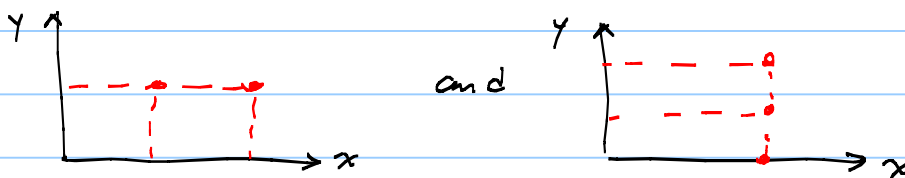
For <u>continuous</u> data : scatterplot     <span style="color:blue">Categ. data, later</span>
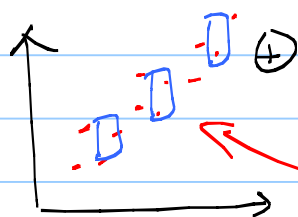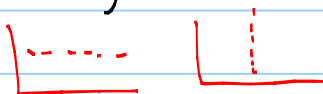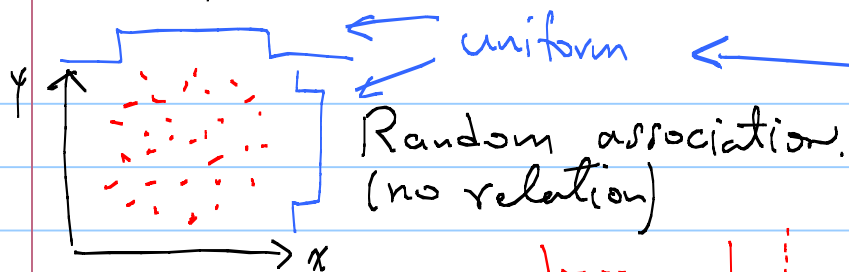


Although one purpose of a scatterplot is to summarize and display the relationship between 2 cont. variables, there is nothing that can fully replace it.

I.e. Given data on 2 vars., do the scatterplot! Of course, histogram each one, too.
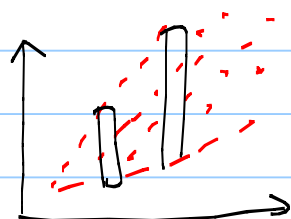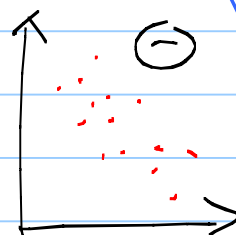


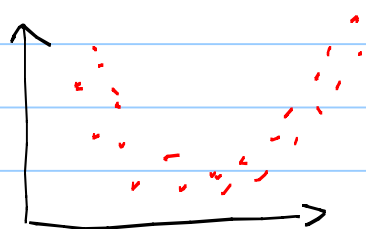Not unusual. In fact, they are common, (and even necessary)

# Scatterplot Museum:

← uniform ←

Random association.
(no relation)

linear, constant variance

Y generally increases with x, but y's variance does not.

$\oplus$

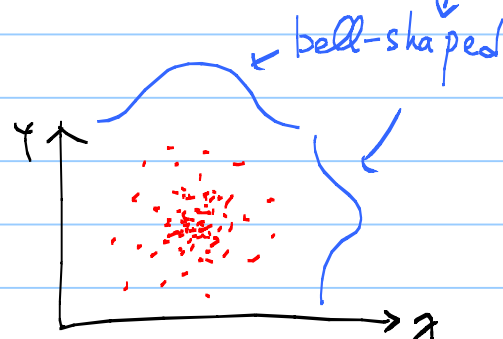$\ominus$

linear, non-constant variance
var. of y changes with x.

generally

non linear (y decreases with increasing x, but only up to some point. Then reverse.)

periodic x & y.

A scatterplot is "the best" device for displaying and studying the relationship (or association) between data on 2 continuous variables.

Q1.) In the scatterplot shown here, there is ☐ relationship between x & y

A) No, B) some C) One cannot say.

The diff. with above is the hist of x (and y).

← bell-shaped

<u>Q</u> Now, how do we quantify the <u>strength</u> of the association between 2 continuous variables?

<u>A</u> there are many measures of strength (like there are different measures of spread of a histogram), and each one captures a different facet of strength.

One popular measure is <u>Pearson's correlation coeff.</u>

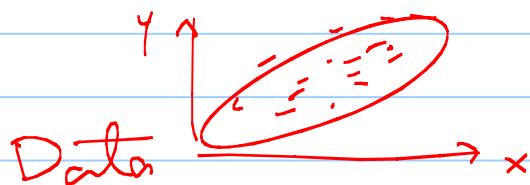denoted $r$ (for sample) and $\rho$ (for distribution)

ie. population

$r$ gives a point estimate of $\rho$.

(like $\bar{x}$ gives a point estimate of $\mu_x = E[x]$

Sample mean

population mean

<u>Q</u> How do we compute it?



Data

<u>A</u>

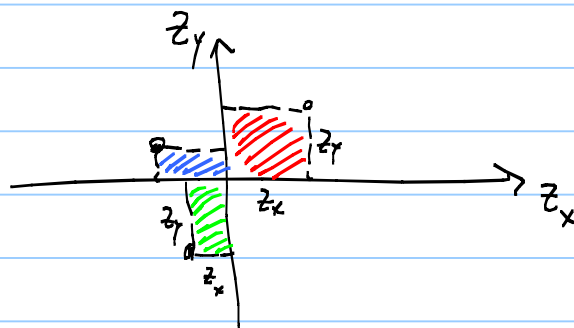| X | Y | | $z_x = \frac{x_i - \bar{x}}{s_x}$ | $z_y = \frac{y_i - \bar{y}}{s_y}$ | $z_x z_y$ |
|---|---|---|---|---|---|
| $x_1$ | $y_1$ | | | | |
| $x_2$ | $y_2$ | | | | |
| $\vdots$ | $\vdots$ | | | | |
| $x_n$ | $y_n$ | | | | |

$\bar{x}, s_x$   $\bar{y}, s_y$

$\frac{1}{n-1} \sum = r$ . "funny Average"

$$\boxed{r_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)}$$   $-1 \leq r_{xy} \leq +1$

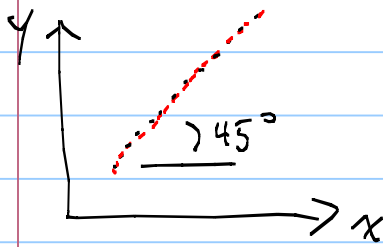$\underbrace{\phantom{xxx}}_{z_x}$  $\underbrace{\phantom{xxx}}_{z_y}$

later!

= Average of "areas" ⟶
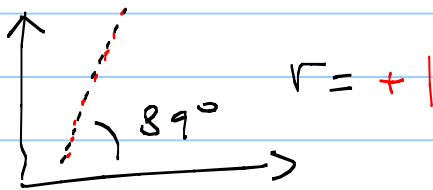


Important: The specific measure of strength that r measures is the "skinniness" of the scatterplot.

generally { fat scatterplot ⟶ r ~ 0
{ skinny " ⟶ r ~ ±1 , But ⟶

r museum:



$r = +1$



$r = +1$

$r = -1$



$r = 0$

$\begin{bmatrix} \text{this involves} \\ \text{some limits} \end{bmatrix}$

$\llcorner \vdots \; r = 0$

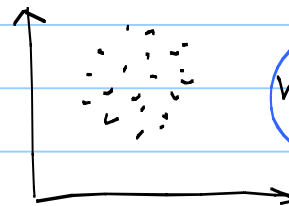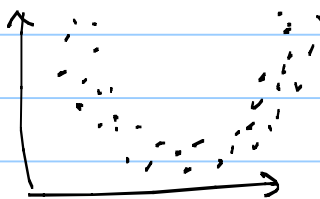$r \sim 0.7, 0.8$

$r \sim -0.7$
$\sim -0.6$

$r \sim 0$

$r \sim 0$

**Important:** r is a summary measure of a scatterplot. As such, some info is lost when you look only at r. Look at The scatterplot (too)!

**hw-lect12-1** Make a scatterplot of The 2 continuous vars in hw-lect1. (By R, or by hand). Describe The relationship.
*If it can't be done, see me!*

**hw-lect12-2** I gave you a formula That defines r. The book gives two others on p. 108.

a) Start from The formula I "derived" in class, and show That it is equal to

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2}\sqrt{\sum (y_i - \bar{y})^2}} \qquad \boxed{I}$$

b) Start from $\boxed{I}$, and show That it is equal to $r = \dfrac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}$, where $S_{xx}, S_{yy}, S_{xy}$ are defined on page 108.

**hw-lect12-3**

Suppose n cases of data on x and y fall exactly on The line $y = mx + b$. Compute The value of r.

Hint: In any of The formulas for r, eliminate all y in favor of x.

**hw-lect12-4**

The $z$'s appearing in The formula for r have two nice properties: Their sample mean is zero, and Their sample variance is 1. prove These!

I.e. show $\bar{z} = \dfrac{1}{n}\sum_i z_i = 0$, $\dfrac{1}{n-1}\sum_i^n (z_i - \bar{z})^2 = 1$