# 11 Regression

- The Correlation Coefficient

- The Least-Squares Regression Line

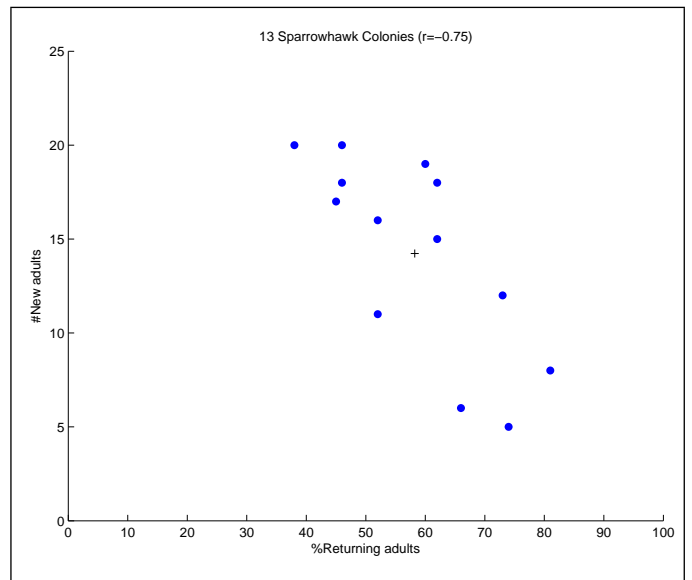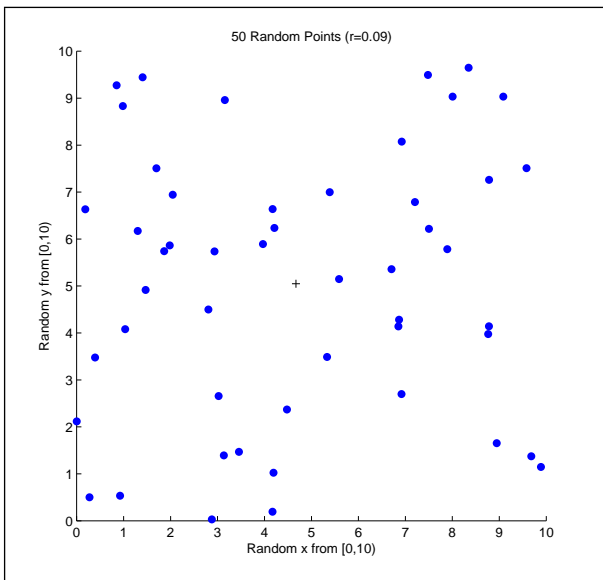## The Correlation Coefficient

### Introduction

A *bivariate* data set consists of $n$ _____,   $(x_1, y_1), \cdots, (x_n, y_n)$.

A *scatterplot* is a _____ of a bivariate data set.

e.g. Here are data for 13 sparrowhawk colonies relating the % of adult sparrowhawks in a colony that return from the previous year and the number of new adults that join the colony:

| %Returning adults | 74 | 66 | 81 | 52 | 73 | 62 | 52 | 45 | 62 | 46 | 60 | 46 | 38 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| #New adults | 5 | 6 | 8 | 11 | 12 | 15 | 16 | 17 | 18 | 18 | 19 | 20 | 20 |

The right-hand scatterplot, below, is from these data. It shows $\cdots$

## The Correlation Coefficient

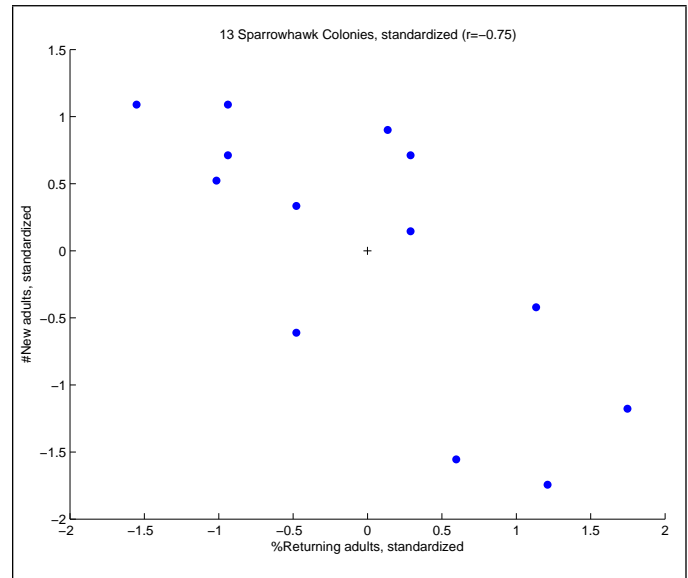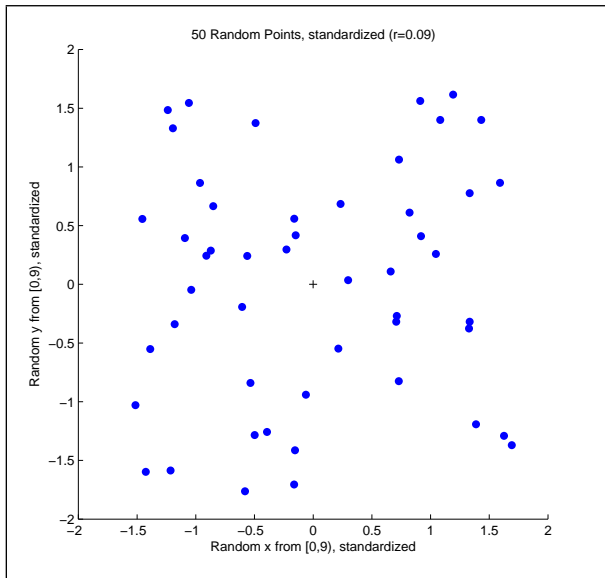The *correlation coefficient*, $r$, measures the _____ and _____ of the linear relationship (if any) between $x$ and $y$:

$$r = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

## An Informal Explanation of $r$

- Start with a scatterplot.

- Shift reference point to _____ by subtracting $\bar{x}$ from each $x_i$ and $\bar{y}$ from each $y_i$.

- Rescale the $x$-axis by dividing each $x$ coordinate by _____,   and rescale the $y$-axis by dividing each $y$ coordinate by $s_y$.

  Now $x$ coordinates, $\frac{x_i - \bar{x}}{s_x}$, have mean _____ and standard deviation _____. $y$ coordinates, $\frac{y_i - \bar{y}}{s_y}$, have the same mean and standard deviation.

- Analyze the sign of the $i^{th}$ term in the sum above, $\left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$, by quadrant:



e.g. For the sparrowhawk data, $r =$ _____.   For the random data, $r =$ _____.

2

**Properties of $r$**

- $-1 \leq r \leq 1$, and

  $r = \pm 1 \implies$ data are _____; $r \approx \pm 1 \implies$ data are _____

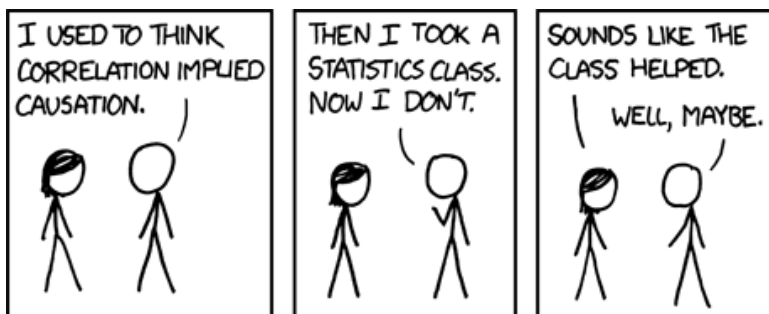  $r \not\approx 0 \implies$ some linear relationship: $x$ and $y$ are *correlated*

  $r > 0 \implies$ slope of line is _____

  $r < 0 \implies$ slope of line is _____

  $r \approx 0 \implies$ no linear relationship: $x$ and $y$ are _____

- $r$ doesn't distinguish between _____ and _____

- $r$ doesn't depend on _____ or _____



I USED TO THINK CORRELATION IMPLIED CAUSATION.

THEN I TOOK A STATISTICS CLASS. NOW I DON'T.

SOUNDS LIKE THE CLASS HELPED. WELL, MAYBE.

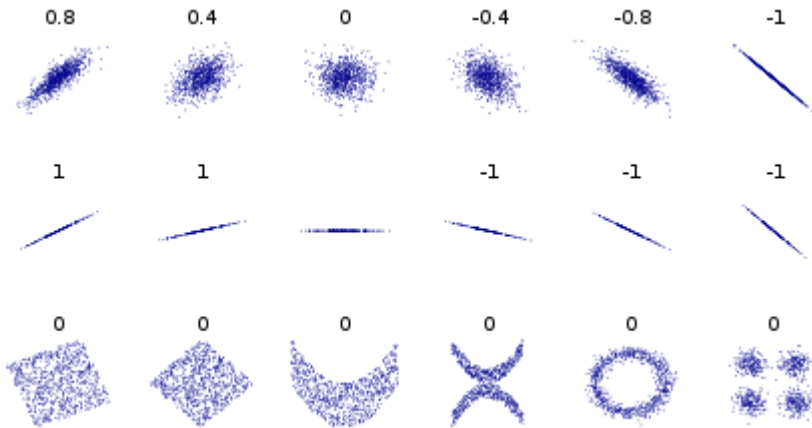http://imgs.xkcd.com/comics/correlation.png

## Cautions

- $r$ measures strength of a *linear* relationship; check scatterplot to avoid using $r$ for a _____

  e.g. The data { (-2, 4), (-1, 1), (0, 0), (1, 1), (2, 4) } fit _____, but $r = 0$ because the data have no _____ relationship (draw).

  e.g. (from http://en.wikipedia.org/wiki/Pearson_product-moment_correlation_coefficient)

  

- $r$ is not resistant to the influence of _____: don't use it for a data set with _____

  e.g. Adding $(0,0)$ to the sparrowhawk data changes $r$ to _____.

- Correlation does not imply causation:

  A _____ (or *lurking*) *variable* is one _____ under consideration that correlates with both the independent and dependent variables of interest.

  e.g.

  – Increasing ice cream sales are correlated with increasing _____ rates. Does ice cream cause _____? _____
  The confounding variable is _____.
  – Sleeping with shoes on is correlated with _____.
  Does sleeping with shoes on cause _____? _____
  The confounding variable is _____.

  If either the independent variable under study, or a _____ confounding variable, affects the dependent variable, then both will seem to by the (_____) criterion of correlation.

  ____ cartoon

# The Least-Squares Regression Line

A _____ line is one that describes how a dependent variable, $y$, changes as an independent variable, $x$, changes in a data set $(x_1, y_1), \cdots, (x_n, y_n)$. We use it to predict $y$ for a given $x$.

The *least-squares regression line* is the line that _____ the data (according to a reasonable criterion).

Notation includes:

- $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$: an unknown true (model) regression line,
  where $\beta_0$ is the $y$-intercept, $\beta_1$ is the slope,
  and $\epsilon_i$ is the $i$th random error

- $y = \hat{\beta}_0 + \hat{\beta}_1 x$: estimated regression line, where

  - $x$: _____ variable

  - $y$: dependent variable

  - $\hat{\beta}_0$: estimated $y$-intercept

  - $\hat{\beta}_1$: estimated _____

- $(x_i, y_i)$: $i^{th}$ data point

- $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$: _____ value of $y$ given $x = x_i$:

- $e_i = y_i - \hat{y}_i$: *residual*, the difference between observed $y_i$ and predicted $\hat{y}_i$; estimates $\epsilon_i$

We predict $y$ from $x$, so minimize vertical error in the "least squares" sense by minimizing a "sum of squared errors"

$$SSE = \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

(Alas, really it should be called a "sum of squared _____.")    Ten lines of calculus gives:

---

For the data set $(x_1, y_1), \cdots, (x_n, y_n)$, the least-squares line is $y = \hat{\beta}_0 + \hat{\beta}_1 x$, where

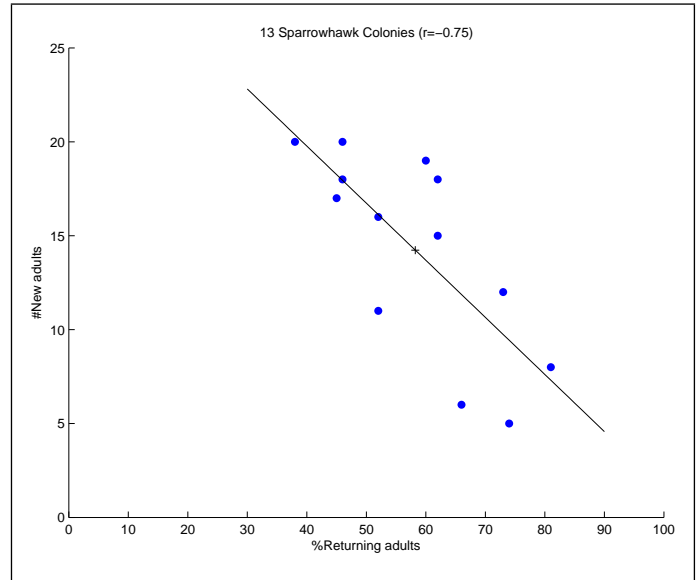$$\hat{\beta}_1 = \frac{s_y}{s_x} r \text{ (slope)}$$
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \text{ ($y$-intercept)}$$

---

e.g. Here again are data for 13 sparrowhawk colonies relating the % of adults in a colony that return from the previous year and the number of new adults that join the colony:

| $x = $ %Returning adults | 74 | 66 | 81 | 52 | 73 | 62 | 52 | 45 | 62 | 46 | 60 | 46 | 38 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $y = $ #New adults | 5 | 6 | 8 | 11 | 12 | 15 | 16 | 17 | 18 | 18 | 19 | 20 | 20 |

Use a calculator to find the least-squares line (recall slope $\hat{\beta}_1 = \frac{s_y}{s_x} r$, $y$-intercept $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$):

- $\bar{x} =$
- $\bar{y} =$
- $s_x =$
- $s_y =$
- $r =$

$\Longrightarrow$

- $\hat{\beta}_1 =$
- $\hat{\beta}_0 =$



13 Sparrowhawk Colonies (r=−0.75)

So our model is $y =$

Or we can do it more directly. (Figure out your _____ labels.)

e.g. Predict the number of new adults in a colony to which 60% of last year's adults return.

$\hat{y} =$_____

## R code for correlation and regression

```
returning = c(74,66,81,52,73,62,52,45,62,46,60,46,38)
new = c( 5,6,8,11,12,15,16,17,18,18,19,20,20)
cor(x=returning, y=new)              # cor() gives correlation
model = lm(new ~ returning)          # lm(y ~ x) gives linear model
model
plot(x=returning, y=new, xlim=c(0, 85), ylim=c(0, 35)) # scatterplot
abline(model)                        # abline() adds line
summary(model)                       # test H_0: beta_i = 0
confint(model)                       # CIs for beta_i
```