

# Lecture 20: Covariance, correlation, and linear regression

Covariance,  
correlation,  
and linear  
regression

## Announcements:

- Reading: Chapter 7 in Vasilj.

# Lecture 20: Covariance, correlation, and linear regression

Covariance,  
correlation,  
and linear  
regression

## Announcements:

- Reading: Chapter 7 in Vasilj.
- Problem set 11 due today.

# Lecture 20: Covariance, correlation, and linear regression

Covariance,  
correlation,  
and linear  
regression

## Announcements:

- Reading: Chapter 7 in Vasilj.
- Problem set 11 due today.
- Today: Covariance, correlation, and linear regression

# Problem. Two continuous variables: are they related?

Covariance,  
correlation,  
and linear  
regression

## Examples

- Parents versus offspring.

# Problem. Two continuous variables: are they related?

Covariance,  
correlation,  
and linear  
regression

## Examples

- Parents versus offspring.
- Growth curves (organs, organisms, populations).

# Problem. Two continuous variables: are they related?

Covariance,  
correlation,  
and linear  
regression

## Examples

- Parents versus offspring.
- Growth curves (organs, organisms, populations).
- Allometric relationships.

# Problem. Two continuous variables: are they related?

Covariance,  
correlation,  
and linear  
regression

## Examples

- Parents versus offspring.
- Growth curves (organs, organisms, populations).
- Allometric relationships.
- Physiological relationships.

# Problem. Two continuous variables: are they related?

Covariance,  
correlation,  
and linear  
regression

## Examples

- Parents versus offspring.
- Growth curves (organs, organisms, populations).
- Allometric relationships.
- Physiological relationships.
- Ecological relationships.



# Correlation versus causation

## Storks



Covariance,  
correlation,  
and linear  
regression

# Correlation versus causation

Covariance,  
correlation,  
and linear  
regression



Babies

# Correlation versus causation

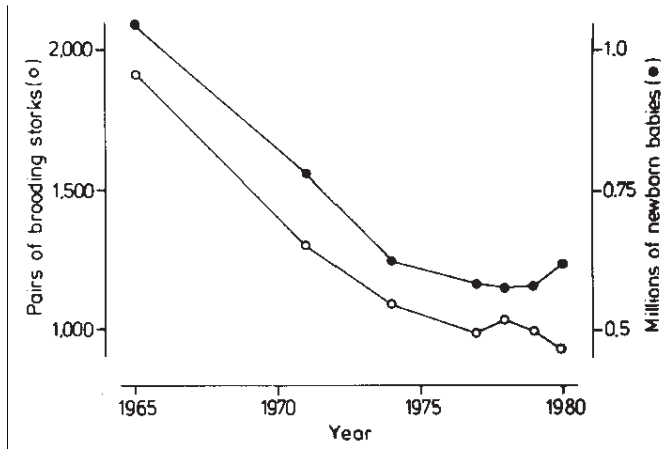
Covariance,  
correlation,  
and linear  
regression



Storks and babies?

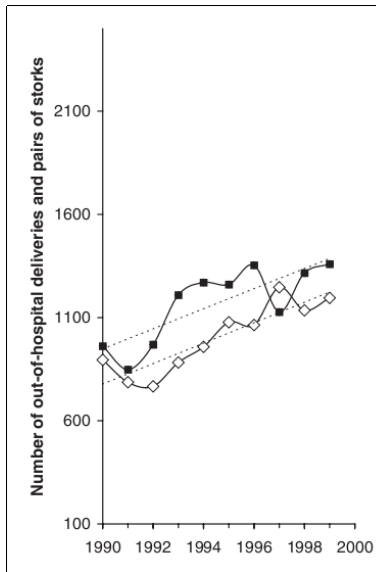
# Correlation versus causation

Covariance,  
correlation,  
and linear  
regression



# Correlation versus causation

Covariance,  
correlation,  
and linear  
regression



# Covariance

Covariance,  
correlation,  
and linear  
regression

$$① \text{Cov}(x, y) = E[(x - \bar{x})(y - \bar{y})]$$

# Covariance

Covariance,  
correlation,  
and linear  
regression

1 
$$\text{Cov}(x, y) = E[(x - \bar{x})(y - \bar{y})]$$

2

$$\text{Cov}(x, y) = \frac{\sum_i [(x_i - \bar{x})(y_i - \bar{y})]}{n - 1}$$

# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- 1  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.



# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- 1  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.
- 2  $\text{Cov}(x, x) = \text{Var}(x)$ , for any variable  $x$ .

# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- 1  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.
- 2  $\text{Cov}(x, x) = \text{Var}(x)$ , for any variable  $x$ .
- 3  $\text{Cov}(x, y) = \text{Cov}(y, x)$ , for any variables  $x$  and  $y$ .

# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- 1  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.
- 2  $\text{Cov}(x, x) = \text{Var}(x)$ , for any variable  $x$ .
- 3  $\text{Cov}(x, y) = \text{Cov}(y, x)$ , for any variables  $x$  and  $y$ .
- 4  $\text{Cov}(x, bx) = b\text{Var}(x)$ , for any variables  $x$  and  $y$  and constant  $b$ .

# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- 1  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.
- 2  $\text{Cov}(x, x) = \text{Var}(x)$ , for any variable  $x$ .
- 3  $\text{Cov}(x, y) = \text{Cov}(y, x)$ , for any variables  $x$  and  $y$ .
- 4  $\text{Cov}(x, bx) = b\text{Var}(x)$ , for any variables  $x$  and  $y$  and constant  $b$ .
- 5  $\text{Cov}(ax, by) = ab\text{Cov}(x, y)$ , for any constants  $a$  and  $b$ .

# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- 1  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.
- 2  $\text{Cov}(x, x) = \text{Var}(x)$ , for any variable  $x$ .
- 3  $\text{Cov}(x, y) = \text{Cov}(y, x)$ , for any variables  $x$  and  $y$ .
- 4  $\text{Cov}(x, bx) = b\text{Var}(x)$ , for any variables  $x$  and  $y$  and constant  $b$ .
- 5  $\text{Cov}(ax, by) = ab\text{Cov}(x, y)$ , for any constants  $a$  and  $b$ .
- 6  $\text{Cov}(x + a, y + b) = \text{Cov}(x, y)$ , for any constants  $a$  and  $b$ .

# Properties of the covariance

Covariance,  
correlation,  
and linear  
regression

- ①  $\text{Cov}(x, a) = 0$ , where  $a$  is any constant.
- ②  $\text{Cov}(x, x) = \text{Var}(x)$ , for any variable  $x$ .
- ③  $\text{Cov}(x, y) = \text{Cov}(y, x)$ , for any variables  $x$  and  $y$ .
- ④  $\text{Cov}(x, bx) = b\text{Var}(x)$ , for any variables  $x$  and  $y$  and constant  $b$ .
- ⑤  $\text{Cov}(ax, by) = ab\text{Cov}(x, y)$ , for any constants  $a$  and  $b$ .
- ⑥  $\text{Cov}(x + a, y + b) = \text{Cov}(x, y)$ , for any constants  $a$  and  $b$ .
- ⑦ If  $x$  and  $y$  are independent, then  $\text{Cov}(x, y) = 0$ .

# Calculating the slope and intercept of the best-fit line

Covariance,  
correlation,  
and linear  
regression

The best fit line is the line that minimizes the residual sum of squares.

1

$$\text{Slope} = b = \frac{\text{Cov}(x, y)}{\text{Var}(x)}.$$

# Calculating the slope and intercept of the best-fit line

Covariance,  
correlation,  
and linear  
regression

The best fit line is the line that minimizes the residual sum of squares.

1

$$\text{Slope} = b = \frac{\text{Cov}(x, y)}{\text{Var}(x)}.$$

2

$$\text{Intercept} = a = \bar{y} - b\bar{x}.$$



# Correlation versus covariance

Covariance,  
correlation,  
and linear  
regression

Correlation coefficient is the scaled covariance.

- 1 The maximum covariance between  $x$  and  $y$  is  $s_x s_y$ .

# Correlation versus covariance

Covariance,  
correlation,  
and linear  
regression

Correlation coefficient is the scaled covariance.

① The maximum covariance between  $x$  and  $y$  is  $s_x s_y$ .

②

$$\text{Cor}(x, y) = r = \frac{\text{Cov}(x, y)}{s_x s_y}.$$

# Correlation versus covariance

Correlation coefficient is the scaled covariance.

- 1 The maximum covariance between  $x$  and  $y$  is  $s_x s_y$ .

- 2

$$\text{Cor}(x, y) = r = \frac{\text{Cov}(x, y)}{s_x s_y}.$$

- 3 The correlation coefficient lies between  $-1$  and  $+1$ .

# Strength of the correlation coefficient

Covariance,  
correlation,  
and linear  
regression

korelacijski koeficijent ( $r$ )	jačina korelacije
0.00 – 0.10	nema
0.10 – 0.25	vrlo slaba
0.25 – 0.40	slaba
0.40 – 0.50	srednja
0.50 – 0.75	jaka
0.75 – 0.90	vrlo jaka
0.90 – 1.00	potpuna

# Simple linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 Calculate the best-fit straight line  $\hat{y} = a + bx$ .

# Simple linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 Calculate the best-fit straight line  $\hat{y} = a + bx$ .
- 2 Calculate the residual sum of squares.  $SS_{res} = \sum (y - \hat{y})^2$ .

# Simple linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 Calculate the best-fit straight line  $\hat{y} = a + bx$ .
- 2 Calculate the residual sum of squares.  $SS_{res} = \sum (y - \hat{y})^2$ .
- 3 Calculate the line sum of squares  $SS_x = \sum (\hat{y} - \bar{y})^2$ .

# Simple linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 Calculate the best-fit straight line  $\hat{y} = a + bx$ .
- 2 Calculate the residual sum of squares.  $SS_{res} = \sum (y - \hat{y})^2$ .
- 3 Calculate the line sum of squares  $SS_x = \sum (\hat{y} - \bar{y})^2$ .
- 4 Calculate the corresponding residual and line mean squares:  $MS_{res} = SS_{res}/df_{res}$ ; and  $MS_x = SS_x/df_x$   
 $df_{res} = n - 2$ ;  $df_x = 1$



# Simple linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 Calculate the best-fit straight line  $\hat{y} = a + bx$ .
- 2 Calculate the residual sum of squares.  $SS_{res} = \sum (y - \hat{y})^2$ .
- 3 Calculate the line sum of squares  $SS_x = \sum (\hat{y} - \bar{y})^2$ .
- 4 Calculate the corresponding residual and line mean squares:  $MS_{res} = SS_{res}/df_{res}$ ; and  $MS_x = SS_x/df_x$   
 $df_{res} = n - 2$ ;  $df_x = 1$
- 5 Form the  $F$  ratio and calculate its tail probability  
 $F = MS_{res}/MS_x$ .

# Simple linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 Calculate the best-fit straight line  $\hat{y} = a + bx$ .
- 2 Calculate the residual sum of squares.  $SS_{res} = \sum (y - \hat{y})^2$ .
- 3 Calculate the line sum of squares  $SS_x = \sum (\hat{y} - \bar{y})^2$ .
- 4 Calculate the corresponding residual and line mean squares:  $MS_{res} = SS_{res}/df_{res}$ ; and  $MS_x = SS_x/df_x$   
 $df_{res} = n - 2$ ;  $df_x = 1$
- 5 Form the  $F$  ratio and calculate its tail probability  
 $F = MS_{res}/MS_x$ .
- 6 Or let R do it all with `anova(lm(y ~ x))`.

# Assumptions of linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 The  $x$  values are measured without any error

# Assumptions of linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 The  $x$  values are measured without any error
- 2 The relationship between  $x$  and  $y$  is linear

# Assumptions of linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 The  $x$  values are measured without any error
- 2 The relationship between  $x$  and  $y$  is linear
- 3 The residuals are normally distributed

# Assumptions of linear regression

Covariance,  
correlation,  
and linear  
regression

- 1 The  $x$  values are measured without any error
- 2 The relationship between  $x$  and  $y$  is linear
- 3 The residuals are normally distributed
- 4 The variance is constant (regardless of the value of  $x$ )

# Interpretation of correlation

Covariance,  
correlation,  
and linear  
regression

