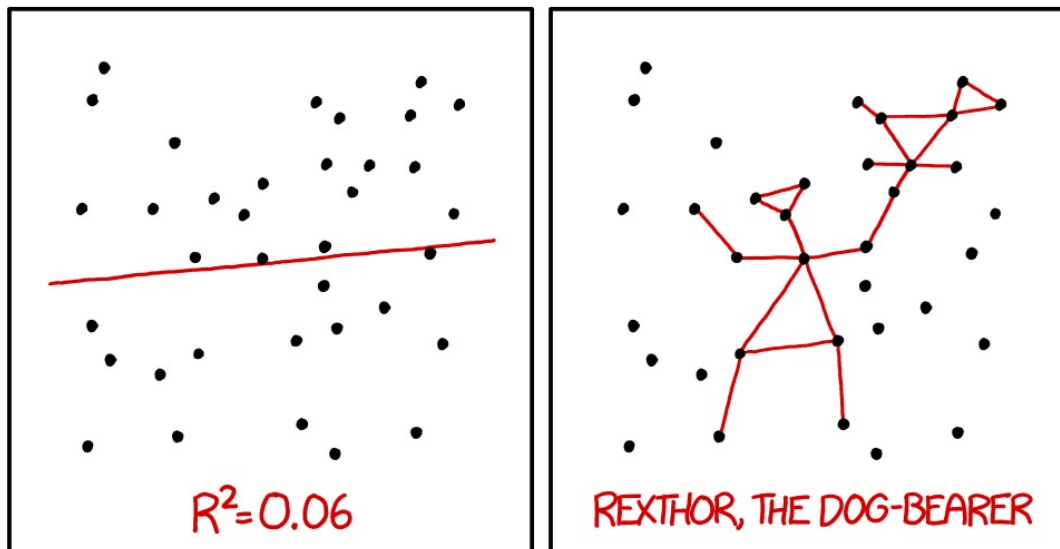


# Stat 88: Probability & Mathematical Statistics in Data Science



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER  
TO GUESS THE DIRECTION OF THE CORRELATION FROM THE  
SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

Lecture 38 : 4/26/2021

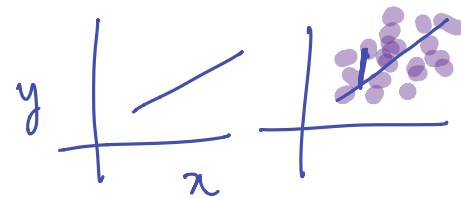
Chapter 11

Correlation

<https://xkcd.com/1725/>

## Equation of the regression line

$$Y = aX + b + \text{error}$$



- $\hat{Y} = \hat{a}X + \hat{b}$

- $\hat{Y}$  is called the fitted value of  $Y$ ,  $\hat{a}$  is the slope,  $\hat{b}$  is the intercept where:

- $\hat{a} = \frac{r\sigma_Y}{\sigma_X}$ ,  $r = E\left[\left(\frac{X-\mu_X}{\sigma_X}\right)\left(\frac{Y-\mu_Y}{\sigma_Y}\right)\right] = E(\underline{Z_X} \times \underline{Z_Y})$

$X^* \quad Y^* \leftarrow$  Text notation

$$Z_X = \frac{X - \mu_X}{\sigma_X}$$

$$r(X, Y)$$

- $\hat{b} = \mu_Y - \hat{a}\mu_X$

$$Z_Y = \frac{Y - \mu_Y}{\sigma_Y}$$

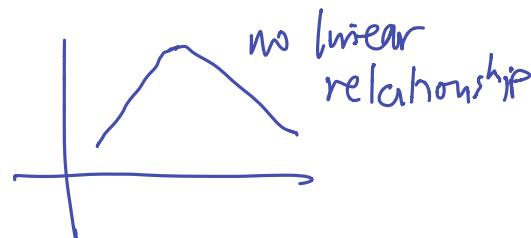
$$D = Y - \hat{Y}$$

# Correlation

$$r = E(Z_X \cdot Z_Y)$$

- The expected product of the deviations of  $X$  and  $Y$ ,  $E(D_X D_Y)$  is called the **covariance** of  $X$  and  $Y$ .
- The problem with using covariance is that the units are multiplied *and* the value depends on the units
- Can get rid of this problem by dividing each deviation by the SD of the corresponding SD, that is, put it in standard units. The resulting quantity is called the **correlation coefficient** of  $X$  and  $Y$ :
- $r(X, Y) = E\left(\frac{D_X D_Y}{\sigma_X \sigma_Y}\right) \Rightarrow E(D_X D_Y) = \text{cov}(X, Y) = r \sigma_X \sigma_Y$
- Note that it is a pure number with no units, and now we will prove that it is always between -1 and 1.

$$E\left(\frac{(X - \mu_X)(Y - \mu_Y)}{\sigma_X \sigma_Y}\right)$$



Bounds on correlation :  $-1 \leq r \leq 1$

$$Z_X = \frac{X - \mu_X}{\sigma_X}$$

- $r = E\left[\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right)\right] = E(Z_X Z_Y)$

$$E(Z_X^2) = \frac{E(X - \mu_X)^2}{\sigma_X^2} = \frac{\sigma_X^2}{\sigma_X^2} = 1$$

- (Note that this implies that  $E(D_X D_Y) = r \sigma_X \sigma_Y$ . We will use this later.)

$$0 \leq E[(Z_X + Z_Y)^2] = E[Z_X^2 + 2Z_X Z_Y + Z_Y^2]$$

$$= E(\underbrace{Z_X^2}_{=1}) + 2E(Z_X Z_Y) + E(\underbrace{Z_Y^2}_{=1})$$

$$0 \leq 1 + 2 \underbrace{E(Z_X Z_Y)}_r + 1$$

$$0 \leq 2 + 2r$$

$$\Rightarrow r \geq -1 \quad (\star)$$

$$(\star) \text{ \& } (\star \star) \text{ imply}$$

$$-1 \leq r \leq 1$$

$$0 \leq E[(Z_X - Z_Y)^2]$$

$$E(Z_X^2 + Z_Y^2 - 2Z_X Z_Y)$$

$$= E(Z_X^2) + E(Z_Y^2) - 2E(Z_X Z_Y)$$

$$= 1 + 1 - 2r$$

$$0 \leq 2 - 2r$$

$$2r \leq 2$$

$$r \leq 1$$

$$(\star \star) \quad 4$$

$$Z_X = \frac{D_X}{\sigma_X} = \frac{X - \mu_X}{\sigma_X}$$

$$D_X = X - \mu_X$$

$$D_Y = Y - \mu_Y$$

$$E(D_X^2) = E((X - \mu_X)^2) = \text{Var}(X) = \sigma_X^2$$

$$E(Z_X^2) = E\left[\left(\frac{D_X}{\sigma_X}\right)^2\right] = \frac{1}{\sigma_X^2} E(D_X^2) = \frac{1}{\sigma_X^2} \cdot \sigma_X^2 = 1$$

Exercise

Suppose  $Y = aX + b$ ,  $a > 0$

$$\mu_Y = E(Y) = a\mu_X + b, \quad \text{Var}(Y) = a^2 \sigma_X^2$$

$$\sigma_Y = a\sigma_X \leftarrow$$

$$r(X, Y) = E(Z_X Z_Y)$$

$$= E\left[\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right)\right]$$

$$= \frac{1}{\sigma_X \cdot a \cdot \sigma_X} E[(X - \mu_X)(aX + b - a\mu_X - b)]$$

$$= \frac{1}{a \sigma_X^2} E[a \cdot (X - \mu_X)^2] = \frac{a \sigma_X^2}{a \cdot \sigma_X^2} = 1$$

$$\begin{aligned} Y &= aX + b \\ \mu_Y &= a\mu_X + b \\ Y - \mu_Y &= aX + b - (a\mu_X + b) \\ &= aX - a\mu_X = a(X - \mu_X) \end{aligned}$$

Exercise : (1) Show that if  $Y = aX + b$ ,  $a < 0$

$$r(X, Y) = -1$$

$$(2) \quad r(aX + b, cY + d) = r(X, Y), \quad ac > 0$$

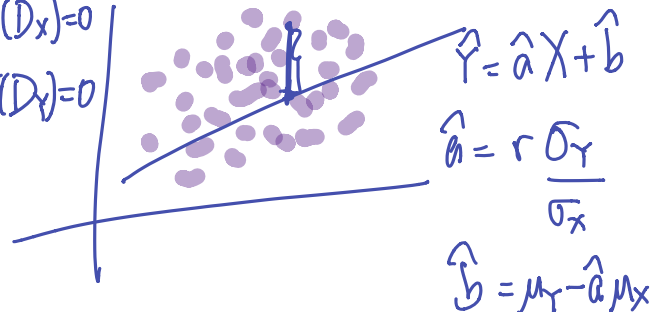
$$-r(X,Y), \text{ as } < 0$$

## Errors in regression

x

$$D_x = X - \mu_x, E(D_x) = 0$$

$$D_y = Y - \mu_y, E(D_y) = 0$$



- The error in regression  $D = Y - \hat{Y}$

$$E(D) = E(Y - \hat{Y})$$

- What is  $E(D)$ ?  $Var(D)$ ?

$$\hat{Y} = \hat{a}X + \hat{b} = \hat{a}X + \mu_y - \hat{a}\mu_x = \hat{a}(X - \mu_x) + \mu_y = \hat{a}D_x + \mu_y$$

$$D = Y - \hat{Y} = Y - \hat{a}D_x - \mu_y = D_y - \hat{a}D_x$$

- Note that we made no assumptions on the distributions of  $X$  &  $Y$ . This means that the residuals average to 0, no matter what the joint distribution of  $X$  &  $Y$ .

- What does the expectation of the error being 0 imply for the residuals?

$$\begin{aligned} E(D) &= E(D_y - \hat{a}D_x) = E(D_y) - \hat{a}E(D_x) = 0 = \mu_D \\ Var(D) &= E((D - \mu_D)^2) = E((D - 0)^2) = E(D^2) \\ Var(D) &= E(D^2) \\ E((Y - \hat{Y})^2) &= \text{Mean squared error for regression} \end{aligned}$$

$$D^2 = \text{MSE for regression} = \text{Var}(D)$$

$$D = D_Y - \hat{a} D_X$$

$$\begin{aligned} E(D^2) &= E[(D_Y - \hat{a} D_X)^2] = E[D_Y^2 - 2\hat{a} D_X D_Y + \hat{a}^2 D_X^2] \\ &= \underbrace{E[D_Y^2]}_{\sigma_Y^2} - 2\hat{a} \underbrace{E(D_X D_Y)}_{r\sigma_X\sigma_Y} + \hat{a}^2 \underbrace{D_X^2}_{\sigma_X^2} \quad \hat{a} = r \frac{\sigma_Y}{\sigma_X} \end{aligned}$$

$$\begin{aligned} E(D^2) &= \sigma_Y^2 - 2r\cancel{\sigma_Y} \cdot r\cancel{\sigma_X} \cdot \cancel{\sigma_Y} + r^2\sigma_Y^2 \cdot \cancel{\sigma_X^2} \\ &= \sigma_Y^2 - 2r^2\sigma_Y^2 + r^2\sigma_Y^2 = \sigma_Y^2 - r^2\sigma_Y^2 \end{aligned}$$

$$\text{Var}(D) = E(D^2) = \text{MSE for regression} = (1-r^2)\sigma_Y^2$$

$$\boxed{SD(D) = \sqrt{1-r^2} \cdot \sigma_Y}$$

4/26/21

$D = Y - \hat{Y}$  = error in regression estimate is called the residual.  
spread in residuals is smaller than spread in  $Y$ .

## Correlation as a measure of linear association

- $D = Y - \hat{Y}$ ,  $E(D) = 0$ ,  $Var(D) = (1 - r^2)\sigma_Y^2$ ,  $SD(D) = \sqrt{1 - r^2} \cdot \sigma_Y$  ( $SD(D) \leq \sigma_Y$ )
- What if the correlation is very close to 1 or -1? What does this tell you about  $X$  &  $Y$ ?  $r \text{ close to } \pm 1 \Rightarrow 1 - r^2 \text{ close to } 0$ ,  $SD(D) \text{ close to } 0$

We know  $E(D) \approx 0$  & if  $SD(D) \text{ close to } 0$

This tells you  $Y$  close to  $\hat{Y} = \hat{a}X + \hat{b} \Rightarrow Y$  is close to being a linear function of  $X$ .

(if  $r = \pm 1$ , then  $Y$  is exactly a linear function of  $X$ )

- What about if the correlation is close to 0? What does this tell you about  $X$  &  $Y$ ?

if  $r \approx 0$ , then  $1 - r^2 \approx 1$

$SD(D) \approx SD(Y)$

then the linear relationship of  $X$  &  $Y$

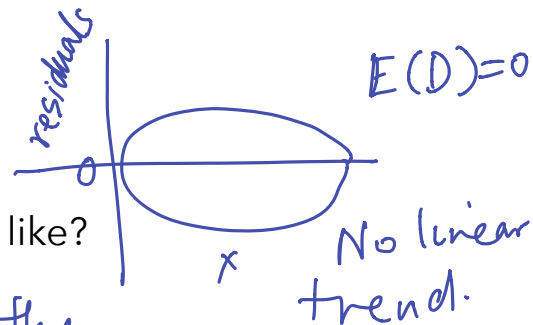
is weak. May as well just use  $\mu_Y$  to predict  $Y$ .



## Residual is uncorrelated with X

- What about  $r(D, X)$ ,  $D = Y - \hat{Y}$ ?
- Intuitively, what should this be? Why?
- What should your residual (diagnostic) plot look like?

$\hat{Y} = \hat{a}X + \hat{b}$      $Y - \hat{Y}$  removes the linear trend with  $X$ .



Exercise : Show that 
$$r(D, X) = \frac{E[(D - \bar{D})(X - \bar{X})]}{\sigma_D \sigma_X} = 0$$

$$E(D \cdot D_X) = ?$$

# The Simple Linear Regression Model

- Regression model from data 8
- Model has two variables: response ( $Y$ ) & ( $x$ ) predictor/covariate/feature variable