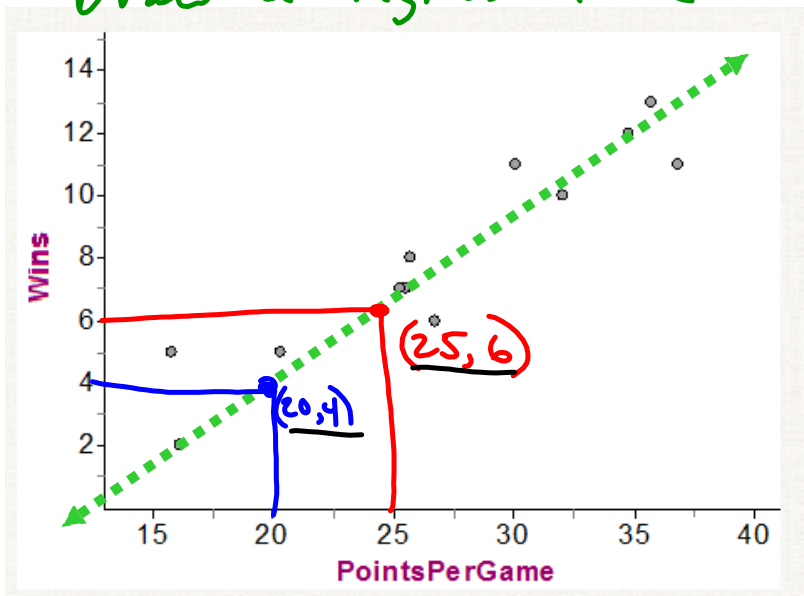


### 3.2 Least Squares Regression

Regression line: line that describes how a response variable,  $y$ , changes as an explanatory variable,  $x$ , changes. It is used to predict  $y$  from  $x$ .



Draw a regression line



A regression line has the form  $\hat{y} = a + bx$  ( $\hat{y} = b_0 + b_1x$ )

- $\hat{y}$  is the predicted value of the response variable for a given explanatory variable,  $x$ .
- $b$  is the slope, which is the amount of change we predict in the response variable for every 1 unit change in the explanatory variable.
- $a$  is the  $y$ -intercept, which is the predicted value of  $y$  when  $x=0$ .

Ex: Write the equation of the regression line for the JEC football example. Interpret the slope &  $y$ -intercept.

$$\hat{y} = -4 + .4x \quad \left( \begin{matrix} 25, 6 \\ x \quad y \end{matrix} \right) \quad b = .4$$

$$\widehat{\text{Win}} = -4 + .4(\text{pointsscore}) \quad \begin{matrix} 6 = a + .4(25) \\ 6 = a + 10 \\ -10 \quad -10 \\ a = -4 \end{matrix}$$

$a$  = When you score zero points, we predict -4 wins.

$b$  = We predict .4 wins for every point per game scored.

Extrapolation is the use of regression for prediction far outside the interval of values of the explanatory variable ( $x$ ) used to obtain the regression line. Beware extrapolating data, often the values will be inaccurate.

Check your understanding, p. 168

1. 40. We predict a rat will gain 40g/week.

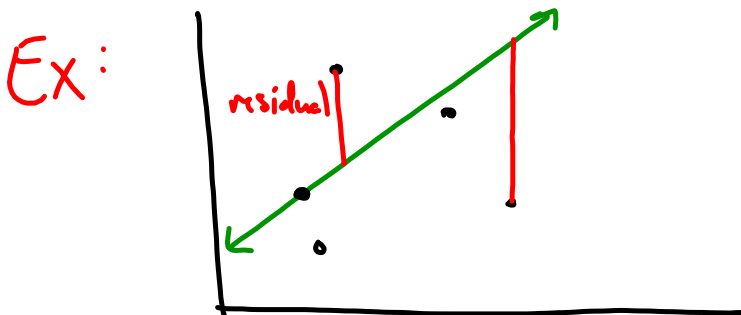
2. 100. At birth, we predict they will weigh 100g.

3.  $\widehat{\text{weight}} = 100 + 40(16)$   
 $= 740 \text{ g.}$

4.  $\widehat{\text{weight}} = 100 + 40(104)$   
 $= 4260 \text{ g}$   
 $\approx 9.4 \text{ lb}$

Residuals: difference in  $y$  between an actual value (point on scatterplot) and a predicted value (point on the regression line).

$$\text{residual} = \text{actual } y - \text{predicted } y = Y - \hat{Y}$$



Ex: Find the residual for 32 points per game (residual) in context.

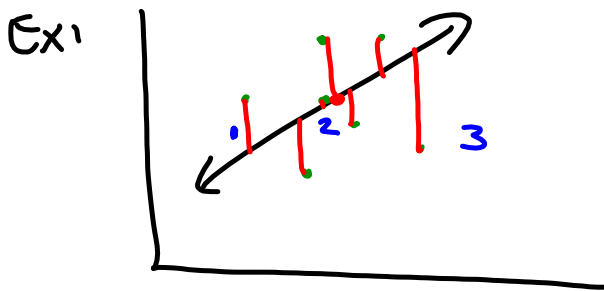
$$\begin{aligned} \text{residual} &= Y - \hat{Y} \\ &= 10 - 8.8 \\ &= 1.2 \end{aligned}$$

$$\begin{aligned} \widehat{\text{Wins}} &= -4 + .4(\text{points per game}) \\ &= -4 + .4(32) \\ &= 8.8 \end{aligned}$$

The predicted number of wins was underestimated by 1.2 wins.

## Least Squares Regression Line (LSRL)

line that makes the squared residuals as small as possible.



### p.170 Activity

- an outlier "pulls" the line towards it.  
The further away from the middle the outlier is, the more effect it will have on the LSRL.
- Always goes through  $(\bar{x}, \bar{y})$

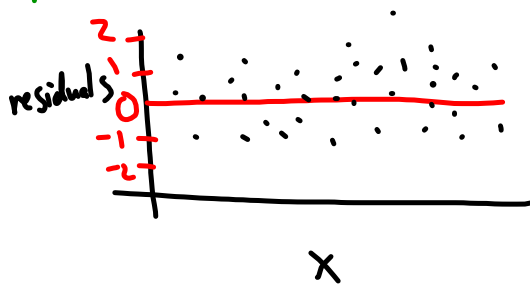
p.172, check your understanding



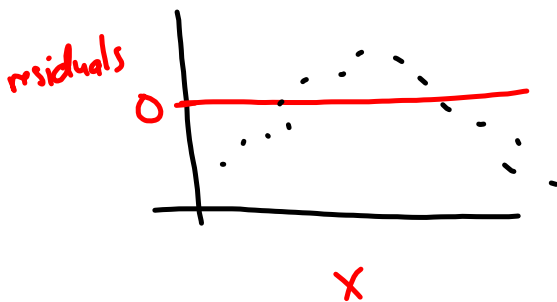
Residual Plot: scatterplot of residuals against the explanatory variable,  $x$ .

This helps us to assess whether a linear model is appropriate.

- A residual plot shows a linear model is appropriate if it is well scattered.



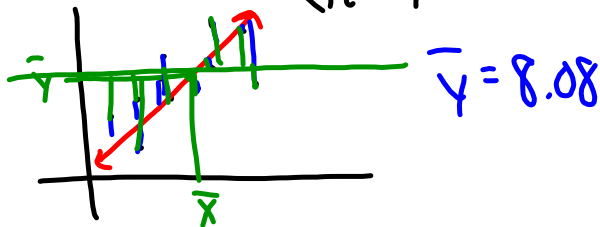
- A residual shows a linear model is not appropriate if it shows a distinct pattern.



Ex: make a residual plot for SEC foot ball data

Coefficient of Determination,  $r^2$ , is the fraction of the variation in the values of  $y$ , accounted for by the LSRL of  $y$  on  $x$ .

$$r^2 = 1 - \frac{\sum (\text{residuals})^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\text{error from LSRL } y \text{ values}}{\text{error from the mean } y \text{ values}}$$



In context,  $r^2$  is interpreted: " — % of the variation in ( $y$ -variable) is accounted for by the linear model relating ( $y$ -variable) to ( $x$ -variable)"

Ex: SEC football example.

88% of the variation in number of wins is accounted for by the linear model relating wins to points per game.

Calculating an LSRL

$$\text{slope} = b = r \cdot \frac{s_y}{s_x}$$

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

## Properties of Correlation + Regression

- ① The distinction of variables as explanatory and response is important regression.
- ② correlation & regression should only be used for straight line relationships.
- ③ correlation and regression are not resistant to outliers. We call an observation an influential point if removing it would drastically change the result of a calculation.
- ④ An association between an explanatory & response variable, even if very strong, is not by itself good evidence changes in  $x$  cause changes in  $y$ .

association  $\neq$  causation