

Correlation and Regression

• In this unit we are exploring bivariate data, which is data involving 2 variables for n observations.

• Specifically we would like to explore relationships between variables. We may represent relationships between variables in several ways. We would like to take a close look at 3:

① - Scatterplots

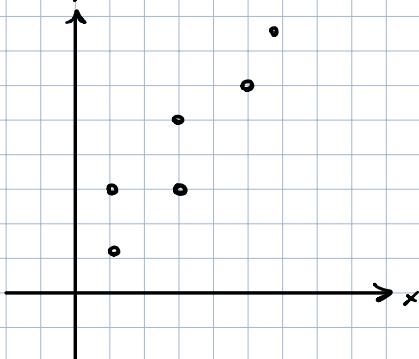


② - Computing a correlation coefficient

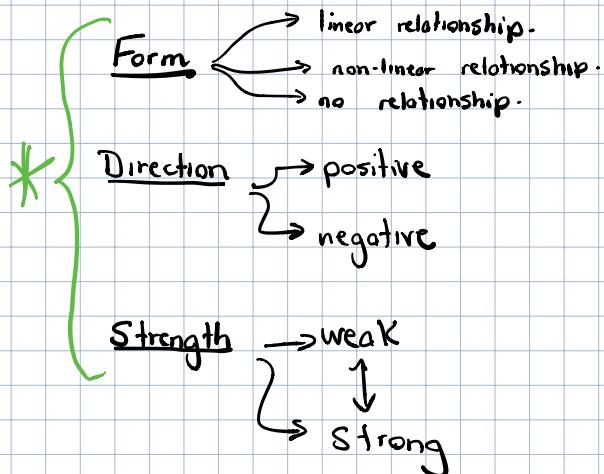
③ - Simple linear Regression (SLR).

Scatter Plots

A scatter plot simply plots one variable on the x -axis, and the other on the y -axis.



We describe a scatterplot in 3-ways:

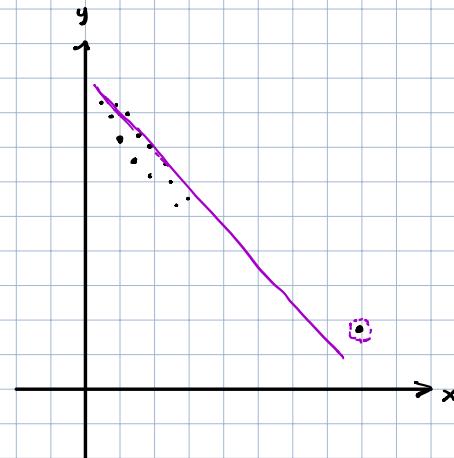
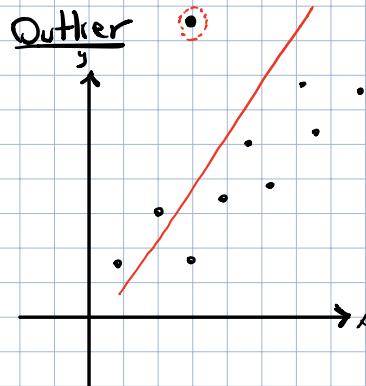


Departures from linearity -

Influential Points → outliers (extreme y-values)

→ points of high leverage (extreme x-values)

Examples

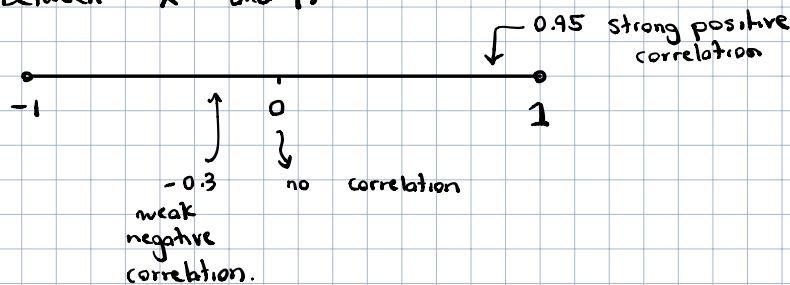


2. Correlation between Variables

- We measure correlation with Pearson's correlation coefficient 'r'!

$$r = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})}{s_x} \frac{(y_i - \bar{y})}{s_y} = \frac{\text{Cov}(X, Y)}{s_x s_y}$$

- r describes both direction and strength of a linear relationship between X and Y.



- correlation is 'unit-less'
- correlation does not imply causation.

3. Simple linear Regression

$$y = a + bx$$

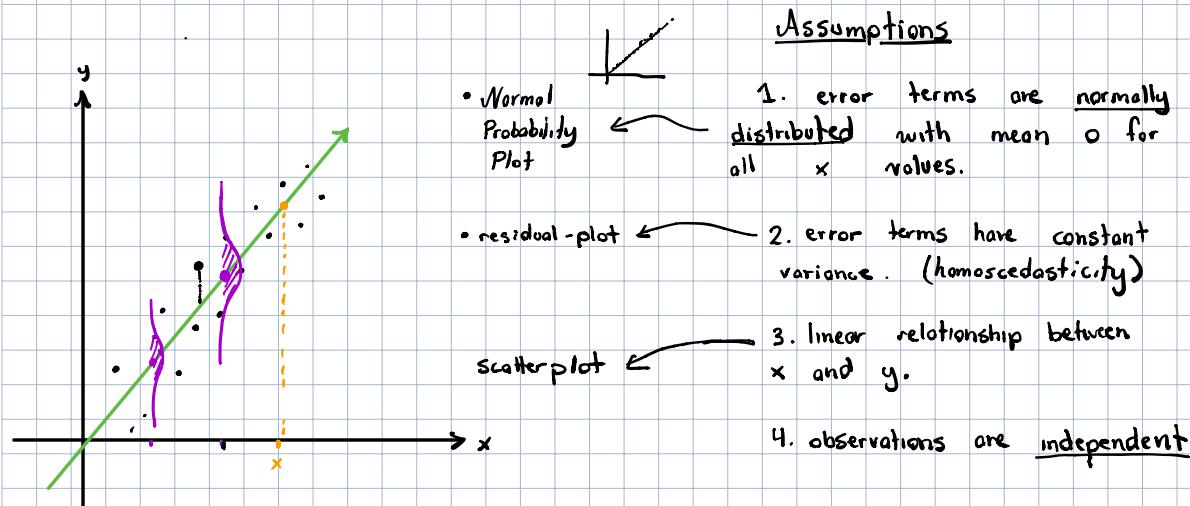
↑ intercept ↑ slope

$b = r \frac{s_y}{s_x}$, $a = \bar{y} - b\bar{x}$

* equation always passes through (\bar{x}, \bar{y}) .

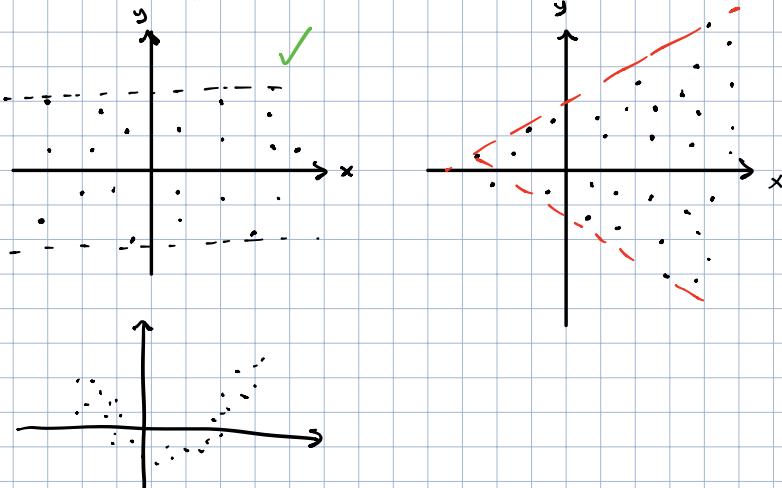
Theoretical Model: $Y = \beta_0 + \beta_1 x + \epsilon$

$$M_{Y|X} = \beta_0 + \beta_1 x$$

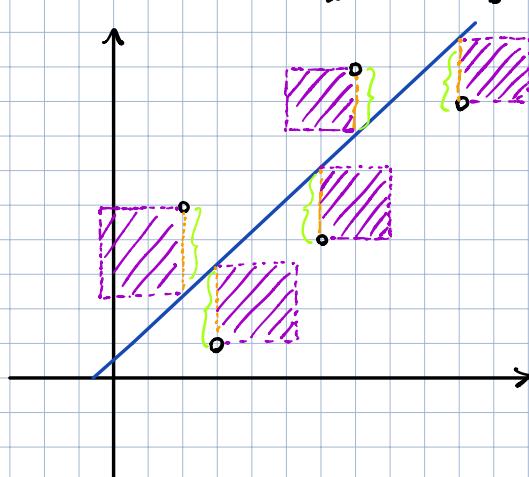


Residual Plots: what is a residual?

$$e_i = y_i - \hat{y}_i$$



- Theoretical Model describes population. We don't actually know true ' β_0 ' and true ' β_1 '. We want to determine best estimate ' $\hat{\beta}_0$ ' and ' $\hat{\beta}_1$ '.



least squares regression
starts with residuals.
take: $e_i = y_i - \hat{y}_i$
true \hat{y} predicted

$$\text{take: } = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x))^2$$

we then minimize this formula with respect to both β_0 and β_1 . (using calculus).

$$\hat{\beta}_1 = r \frac{S_y}{S_x} = \frac{SS_{xy}}{SS_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

How to interpret β_0 and β_1 ?

a/ β_0 - Describes predicted value of y when x is 0.

• Oftn meaningless. Ex: height vs. weight.

$$\text{slope} = \frac{\text{rise}}{\text{run}} \\ = \frac{\text{rise}}{\text{one}}$$

b/ β_1 : Describes the average increase/decrease in y for a 1-unit change in x .

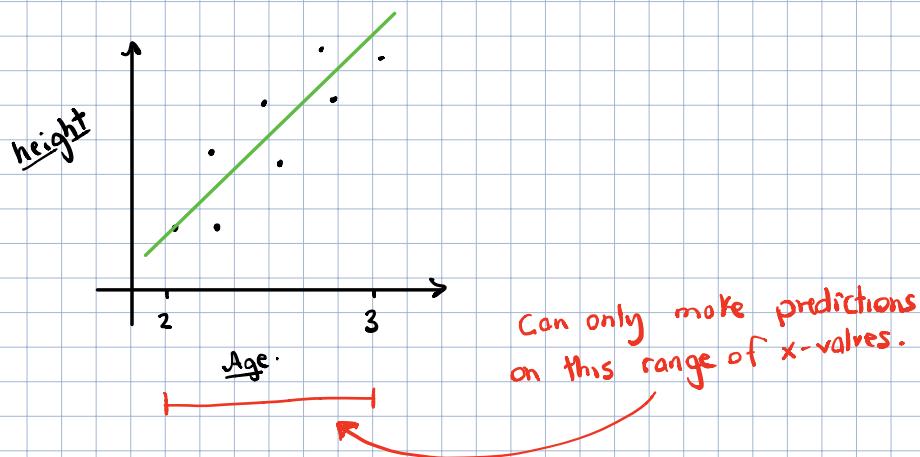
Prediction using a SLR model

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\hat{y} = a + bx$$

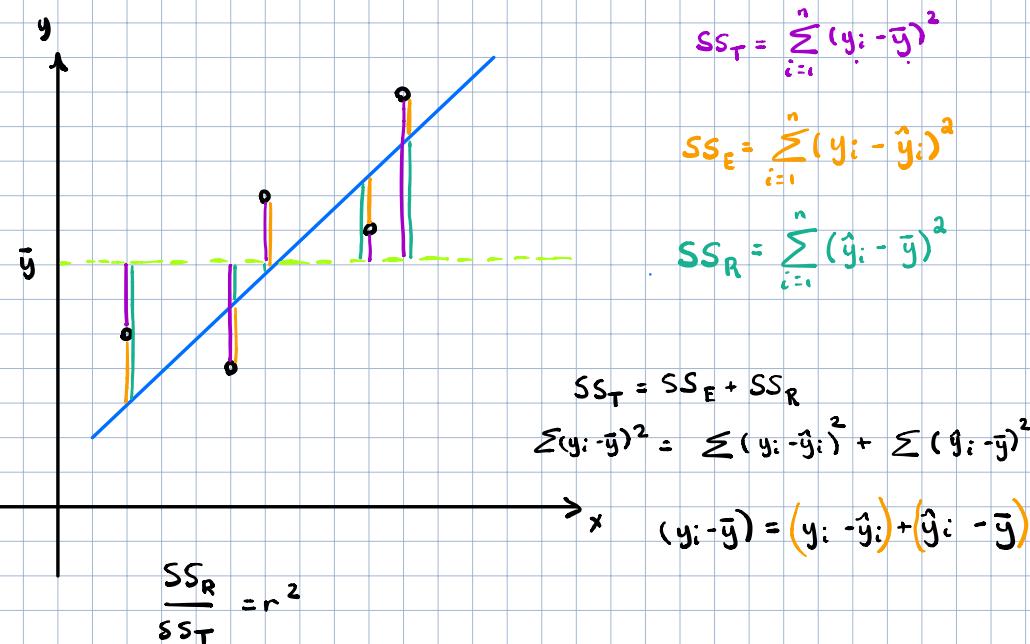
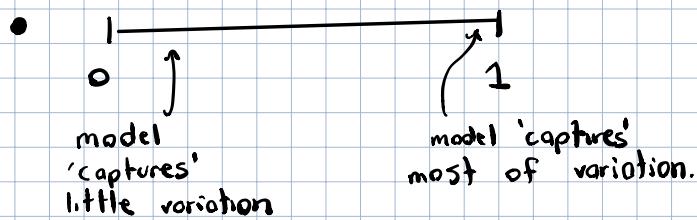
\hat{y} : predicted value of y given x .

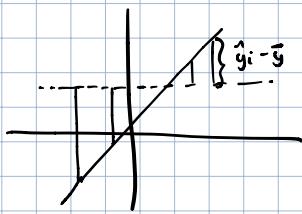
Extrapolation: predicting values that fall well beyond the range of the data.



Coefficient of determination: r^2

- literally correlation coefficient squared.
- r^2 describes the proportion of variation in y (response variable) that can be explained by its linear relationship with the x (explanatory) variable.





$$1 = \frac{SSE}{SST} + \frac{SSR}{SST} \rightarrow r^2$$

$$r^2 = 1 - \frac{SSE}{SST}$$

The Standard error

- describes 'common' error term for each value of x .

Denoted S , S_e , SE . Measure of the amount the regression equation over/under predicts on average.

$$S = \sqrt{\frac{\sum (y_i - \hat{y})^2}{n-2}}$$

Transformations

- we may 'transform' x and y variables to make a relationship more linear. After we apply a transformation there are several ways to check if it worked.

did it work?

- increase randomness in residual plot.
- r^2 may be closer to 1.
- relationship in scatterplot is visually more linear.

• We could keep going... and if we did...

- Notice $\hat{\beta}_0$, and $\hat{\beta}_1$ are statistics created from sample data. They will vary from sample to sample.

→ we would then look at inference for $\hat{\beta}_0$ and $\hat{\beta}_1$

↳ H.T. for $\hat{\beta}_0$ and $\hat{\beta}_1$

↳ C.I. for $\hat{\beta}_0$ and $\hat{\beta}_1$.

In Case you Wondered Least Squares Regression

NOT Mandatory

$$\text{A little backwards: } \hat{\beta}_1 = r \frac{s_y}{s_x} = \left(\frac{1}{n-1} \right) \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y} \cdot \frac{s_y}{s_x} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{s_x^2} \cdot \frac{1}{n-1}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{\sum (x_i - \bar{x})^2}{n-1} \cdot (s_x^2)}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i y_i - \bar{y} x_i - \bar{x} y_i + \bar{x} \bar{y})}{\sum (x_i^2 - 2\bar{x} x_i + \bar{x}^2)}$$

$$= \frac{\sum x_i y_i - \bar{y} \sum x_i - \bar{x} \sum y_i + n \bar{x} \bar{y}}{\sum x_i^2 - 2\bar{x} \sum x_i + n \bar{x}^2}$$

$$= \frac{\sum x_i y_i - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y}}{\sum x_i^2 - 2n \bar{x}^2 + n \bar{x}^2}$$

$$= \boxed{\frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}}$$

this is
what we get
in our proof.

Now let's prove:

$$\text{Sum of least squares: } \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \\ = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

then take partial derivatives to minimize w.r.t β_0, β_1

$$\frac{\partial}{\partial \beta_0} (\sum e_i^2) = -2 \sum y_i - \beta_0 - \beta_1 \sum x_i \\ = -2 \sum y_i + 2n \beta_0 + 2 \beta_1 \sum x_i$$

$$2 \sum y_i = 2n \beta_0 + 2 \beta_1 \sum x_i \\ \text{divide by } 2n$$

$$\boxed{\bar{y} = \beta_0 + \beta_1 \bar{x}} \quad \text{equation 1}$$

$$\frac{\partial}{\partial \beta_0} (\sum e_i) = -2 \sum x_i (y_i - \beta_0 - \beta_1 x_i)$$

$$= -2 \sum x_i y_i + 2 \beta_0 \sum x_i + 2 \beta_1 \sum x_i^2$$

$$2 \sum x_i y_i = 2 \beta_0 \sum x_i + 2 \beta_1 \sum x_i^2$$

$$\sum x_i y_i = \beta_0 \sum x_i + \beta_1 \sum x_i^2$$

$$\sum x_i y_i = (\bar{y} - \beta_1 \bar{x}) \sum x_i + 2 \beta_1 \sum x_i^2$$

$$\sum x_i y_i = \bar{y} \sum x_i - \beta_1 \bar{x} \sum x_i + \beta_1 \sum x_i^2$$

$$\sum x_i y_i = n \bar{x} \bar{y} - \beta_1 n \bar{x}^2 + \beta_1 \sum x_i^2$$

$$\beta_1 \sum x_i^2 - \beta_1 n \bar{x} = \sum x_i y_i - n \bar{x} \bar{y}$$

$$\beta_1 (\sum x_i^2 - n \bar{x}^2) = \sum x_i y_i - n \bar{x} \bar{y}$$

$$\boxed{\beta_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = r \frac{s_y}{s_x}}$$