

Ref: Ch 5, Faraway text

If  $X$  does cause  $Y$  ( $X \rightarrow Y$ ), then we should observe some association (not necessarily linear) between  $X$  &  $Y$ , but the converse is not necessarily true

"correlation does not imply causation"

Theories of causality differ between disciplines, but all share some common features:

- underlying theory: the "science" suggest some mechanism by which  $X$  might cause  $Y$  & also rules out alternative causes (say  $Z$ )



- temporal order:  $X$  must precede  $Y$   
 (so that it is  $X \rightarrow Y$ , not  $Y \rightarrow X$ )
- association:  $X \rightarrow Y$  will usually result in some correlation (linear association) between  $X$  &  $Y$   
 note: relationship may not be linear

If we discover "associations" in observational data & suspect that it is because  $X \rightarrow Y$  then in the next iteration of the research process we might <sup>try</sup> some more structured approach (e.g. a designed experiment)

Coefficient of Determination ( $R^2$ )

(2)

"Proportion of the variation in  $Y$  that can be explained by the model involving the  $X(s)$ "

- In the R output this is Multiple R-squared; called "multiple" as it does generalise to multiple regression of  $Y$  on 2 or more  $X$ s

$$R^2 = \frac{SS_{\text{Regression}}}{SS_{\text{Total}}} = 0.7388 \text{ or } 74\%$$

$$= 1 - \frac{SS_{\text{Error}}}{SS_{\text{Total}}} \leftarrow s^2 = \hat{\sigma}^2 = \sum e_i^2$$

$$= 1 - \frac{SS_{\text{Error}}}{SS_{\text{Total}}} \leftarrow s_y^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2$$

so,  $SS_T = (n-1) * \text{var}(y)$

NB:  $R^2 = (r)^2$  coefficient of correlation between  $Y$  &  $X$

- Adjusted  $R^2$  (adjusted for the degrees of freedom)

$$\bar{R}^2 = 1 - \frac{MS_{\text{Error}}}{MS_{\text{Total}}} = 1 - \frac{SS_{\text{Error}}/df_{\text{Error}}}{SS_{\text{Total}}/df_{\text{Total}}}$$

$$= \dots = R^2 - \frac{(1-R^2) \frac{df_{\text{Regression}}}{df_{\text{Error}}}}{\text{adjustment factor}}$$

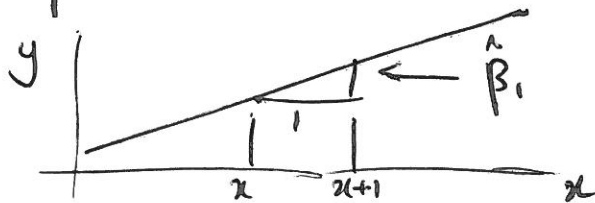
cf.  $F = \frac{MS_{\text{Regression}}}{MS_{\text{Error}}} \sim F_{k, n-p}$  Overall F statistic  
known dist.

These are all summary measures. But the F statistic has some advantages

- like  $\bar{R}^2$  it does adjust for the df (Ex: shown you derive  $\bar{R}^2$  from F)
- F is comparable to a known standard distribution (still have to choose  $\alpha$ )

## Interpreting the regression coefficients

- Interpretation of  $\hat{\beta}_1$  (or any slope coefficient) is "the expected increase in  $Y$  as  $X$  increases by 1"



- Interpretation of  $\hat{\beta}_0$  is "the expected value of  $Y$  when  $X = 0$ " (ie it is the intercept coefficient)

A 95% confidence interval for  $\beta_1$  is

$$\underbrace{\hat{\beta}_1}_{\text{estimate}} \pm \underbrace{t_{\text{error df}}(0.975)}_{\text{critical value}} \cdot \underbrace{\text{se}(\hat{\beta}_1)}_{\substack{\text{Standard error} = \frac{\Delta}{\sqrt{S_{xx}}} \\ \text{(for SLR)}}}$$

Similarly a 95% confidence interval for  $\beta_0$  is

$$\hat{\beta}_0 \pm t_{\text{error df}}(0.975) \cdot \text{se}(\hat{\beta}_0)$$

$\nearrow \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$