

$p=1; y=R$

$n^*(x) = \beta_0 + \beta_1 x, y = n^*(x) + \epsilon$

$g(x) = b_0 + b_1 x$

A: OLS  $\Rightarrow b_0 = \bar{y} - r \frac{s_y}{s_x} \bar{x}, b_1 = r \frac{s_y}{s_x}$

$y = g(x) + \epsilon \leftarrow (\epsilon = \text{OLS error types})$

How well does  $g$  predict?

$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - g(x_i))^2$

$y$  means squared

$MSE = \frac{1}{n-2} SSE$

mean squared error

$RMSE = \sqrt{MSE}$

Root mean squared error

$\epsilon$  is pulled from a normal distribution

$g(x) \pm RMSE$

95% prediction model

\* now evaluate your model is

• IF  $R^2 \uparrow \Leftrightarrow SSE \downarrow \Leftrightarrow MSE \downarrow \Leftrightarrow RMSE \downarrow$   
IF  $RMSE = 0 \Leftrightarrow SSE = 0 \Leftrightarrow MSE = 0 \Leftrightarrow R^2 = 1$   
• IF  $R^2 = 1$ ; RMSE is large

•  $R^2$  vs RMSE  
WHICH IS MORE IMPORTANT?  
RMSE because it shows how off your model is in units which is more useful than % (usually).  
 $\epsilon = RMSE$

$p=1$   
 $\downarrow$   
 $x_1 = \text{Binary } 0, 1$   
 $x \in X = \{\text{red, green}\}$

$x = 1 \Rightarrow \text{green}$

$H = \{w_0 + w_1 x : w_0, w_1 \in \mathbb{R}\}$

$\hat{y} = g(x) = b_0 + b_1 x$

$\bar{x} = \frac{\sum x_i}{n} = \frac{n_{\text{green}}}{n} = \text{proportion of green} \rightarrow p$

$\bar{y} = \frac{\sum y_i}{n} = \frac{\sum_{i \in \text{green}} y_i + \sum_{i \in \text{red}} y_i}{n} = \frac{\sum_{i \in \text{green}} y_i}{n} + \frac{\sum_{i \in \text{red}} y_i}{n}$

$= p \bar{y}_{\text{green}} + (1-p) \bar{y}_{\text{red}}$

$b_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{n \bar{y}_{\text{green}} - n p \bar{y}}{n - n p^2} = \frac{\bar{y}_{\text{green}} - p \bar{y}}{1 - p}$

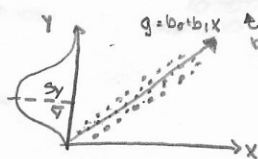
$= \frac{(1-p) \bar{y}_{\text{green}} - (1-p) \bar{y}_{\text{red}}}{1-p}$

$b_0 = \bar{y} - b_1 \bar{x} = (\bar{y}_{\text{green}} - (1-p) \bar{y}_{\text{red}}) - (\bar{y}_{\text{green}} - \bar{y}_{\text{red}}) p + (1-p) \bar{y}_{\text{red}} = \bar{y}_{\text{red}}$

consider the null model

$g(x) = \bar{y}$

$SSE_0 = \sum_{i=1}^n (y_i - \bar{y})^2 = SST = (n-1) s_y^2$   
sum of squares total



Model must beat SST

distr of  $\epsilon$ 's (residuals)  
+ if "thinner", better model

$SSE = \sum_{i=1}^n e_i^2 = (n-1) s_e^2$

$AS^2 = s_y^2 - s_e^2$

reduction in variance / variance explained

$R^2 = \frac{AS^2}{s_y^2} = \frac{s_y^2 - s_e^2}{s_y^2} = 1 - \frac{s_e^2}{s_y^2} = \frac{(n-1) s_y^2 - SSE}{(n-1) s_y^2} = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$

\* OLS will always be non-negative

• can  $R^2 > 1$ ?

$1 - \frac{SSE}{SST} \geq 1 \Rightarrow SSE \leq 0$

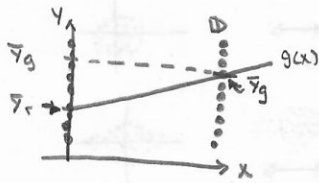
since  $SSE \geq 0$  and  $SST \geq 0$

$\frac{SSE}{SST} \geq 0$

$R^2 = 1$  only when  $SSE = 0$

• can  $R^2 = 0$   
 $SSE = SST$   
 $g = \bar{y}$

• can  $R^2 < 0$   
 $SSE > SST$   
bad model  
 $g$  predicts worse than  $\bar{y}$   
 $g_0 > \bar{y}$



$g(x) = 0.7 + 0.2x$

$\bar{y}_{\text{green}} - \bar{y}_{\text{red}} = \frac{\sum_{i \in \text{green}} y_i}{n - n_g} + \frac{\sum_{i \in \text{red}} y_i}{n - n_g} = \frac{\sum y_i}{n - n_g} = \frac{\bar{y}}{1 - p}$

$L = 3$

$X \in \{red, green, blue\}$

$X_1 = \mathbb{I}_{X=green} \in \{0,1\}$

$X_2 = \mathbb{I}_{X=blue} \in \{0,1\}$

$\hat{y} = b_0 + b_1 X_1 + b_2 X_2$

$\hat{y} = \bar{y}_{red} + (\bar{y}_g - \bar{y}_r) X_1 + (\bar{y}_b - \bar{y}_r) X_2$   
 $\quad \quad \quad -2 \quad \quad \quad -0.4$

$X \in \{low, med, high\}$  ordinal argument

$X_1 = \mathbb{I}_{X=low}$

$X_2 = \mathbb{I}_{X=high}$

$\hat{y} = \begin{cases} \bar{y}_L & \text{if } X=low \\ \bar{y}_M & \text{if } X=medium \\ \bar{y}_H & \text{if } X=high \end{cases}$

what if you want to constrain

$\hat{y}(low) < \hat{y}(medium) < \hat{y}(high)$

A: OLS will not give you  $\hat{y}$  we want monotonicity

consider the r.v's  $X, Y$

they are dependent if (associated if)

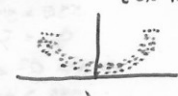
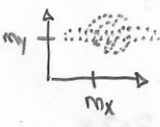
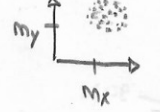
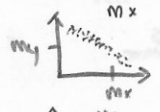
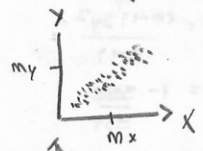
$\exists x_1, x_0 \text{ s.t. } P(Y|X=x_1) \neq P(Y|X=x_0)$

$E_{xy} = cov[X, Y] = \frac{S_{xy}}{S_x S_y}$  estimated by  $r$

$S_{xy} = cov[X, Y] = E[(X - M_X)(Y - M_Y)]$  estimated by  $\sum X_i Y_i$

let  $X_c = X - M_X; Y_c = Y - M_Y$

$S_{xy} = E[(X - M_X)(Y - M_Y)] = E(X_c Y_c) = E(Z)$



$E(Z) > 0$

$E(Z) < 0$

$E(Z) \approx 0$

correlation = "association"  
 $\hookrightarrow$  linear associations

$r^2 = R^2$  for  $p=1$   
 sample correlation coefficient

$E(Z) \approx 0 \Rightarrow \text{correlation} = 0$

$X$  &  $Y$  dependent / associative