

# REGRESSION

HAVE OBSERVATIONS OF A  
DEPENDENT VARIABLE  $Y$ , AND  
 $K$  INDEPENDENT VARIABLES  
 $X_1, X_2, \dots, X_K$

CONSIDER LINEAR RELATIONSHIP

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_K X_K$$

"Y-HAT"

INTERCEPT  
↑  
SCOPE COEFF  
↑  
↓  
PREDICTED, OR FITTED VALUE

EFFECT OF UNIT  
INCREASE IN ONE  
VARIABLE  $X_i$   
HOLDING OTHER  
VARIABLES FIXED

LET  $\hat{y}_i$  BE THE FITTED  
VALUE FOR  $i$ 'TH OBSERVATION

LET  $y_i$  BE OBSERVED VALUE

LET  $e_i = y_i - \hat{y}_i$  BE ERROR  
OR RESIDUAL

USUAL CRITERION:

MINIMIZE  $\sum_{i=1}^n e_i^2$

"LEAST SQUARES" REGRESSION

- ✓ COMPUTATIONALLY INTENSIVE  
[UP TO REASONABLE SIZE]
- ✓ PERMITS INFERENCE  
UNDER APPROPRIATE ASSUMPTIONS



## MEASURE OF FIT

$$SSY = SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

"sum of squares"      ↑      "TOTAL"

$$SS\hat{y} = SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

↑      "REGRESSION"

$$SSE = \sum_{i=1}^n e_i^2$$

↑      "ERROR"

Fact For L-S fit

$$SSE + SSR = SST$$

$$\frac{SSE}{SST} + \frac{SSR}{SST} = 1$$

u  $\downarrow$  "FRACTION OF VARIATION IN Y 'EXPLAINED BY' INDEP. VARIABLES"

$R^2$

\* USING LINEAR RELATIONSHIP

4

ANOTHER MEASURE OF FIT

$$s_e = \sqrt{\frac{SSE}{n - (k+1)}} \dots$$

"STD. ERROR OF ESTIMATE"  
(ESSENTIALLY SAMPLE STD. DEV. OF RESIDUALS)

## INFERENCE

IDEA: WANT MEASURE FOR  
POSSIBLE ERROR IN  
OUTPUT  $b_0, b_1, \dots, b_k$

NEED TO ASSUME A "MODEL"

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \epsilon$$

$\epsilon$  IS NORMAL R.V.

$$E[\epsilon] = 0$$
$$\text{VAR}[\epsilon] = \sigma_\epsilon^2$$

$\epsilon$  IS INDEPENDENT OF  
 $X_1, X_2, \dots, X_k$



5

UNDER MODEL ASSUMPTION,  
OUTPUTS OF L-S REGRESSION  
HAVE KNOWN "SAMPLING  
DISTRIBUTIONS".....

BOTTOM LINE: CONSTRUCT  
CONFIDENCE INTERVALS  
FOR  $\beta_0, \beta_1, \dots, \beta_K$

HOW TO EVALUATE  
REGRESSION MODELS:

(1) INTERPRETATION OF  
EQUATION

(2) RESIDUALS

- PLOTS AGAINST  $X_i$
- HISTOGRAM