

TODAY: LOGISTIC REGRESSION

IDEA: MULTIPLE REGRESSION
WITH DUMMY (0/1)
DEPENDENT VARIABLE

- MANY APPLICATIONS
- REGULAR L-S REGRESSION
DOESN'T MAKE SENSE...

EX: COMMUNITY COLLEGE DATA
(RATE)
CONCERNED WITH STUDENT
SUCCESS
GPA [GPA < 2.0]

EX: FORECAST GPA = 2.3

$$P(\text{GPA} < 2.0) = P\left(\frac{\text{GPA} - E[\text{GPA}]}{\sigma_{\text{GPA}}}\right)$$

STD ERROR
FROM REGRESSION
OUTPUT

$$\frac{2.0 - 2.3}{.7759}$$

$$= P(Z < -\frac{3}{.7759})$$

$$= P(Z < -3.87)$$

~ 75%

[From z table]

NOTE: σ_{GPA} ACTUALLY VARIES WITH VALUES OF INDEP. VARIABLES, BUT FOR TYPICAL VALUES IS DOMINATED BY STD. ERROR S_e

LOGISTIC REGRESSION

IDEA: LET $p = P(Y=1)$
(FOR SET OF VALUES
OF indep VARIABLES)

- INSIGHT OF PREDICTING Y (0/1),
PREDICT P
- TURNS OUT TO BE BETTER
TO WORK WITH
ODDS RATIO

$$OR = \frac{P}{1-P}$$

EX: $p = .5$, $OR = \frac{.5}{.5} = 1$ "EVEN"

$p = .9$ $OR = \frac{.9}{.1} = 9$ "9 TO 1
IN FAVOR"

$p = .2$ $OR = \frac{.2}{.8} = \frac{1}{4}$ "4 TO 1
AGAINST"

$\rightarrow \frac{1}{OR} = \frac{1-P}{P} = \frac{1}{P} - 1$

$\frac{1}{OR} + 1 = \frac{1}{P}$

$\frac{OR+1}{OR} = \frac{1}{P}$

$P = \frac{OR}{1+OR}$

- TURNS OUT EVEN BETTER TO
WORK WITH LN OF OR

$$\underline{\text{LOR}} = \text{LN}(\text{OR}) = \text{LN}\left(\frac{P}{1-P}\right)$$

$$e^{\text{LOR}} = e^{\text{LN}(\text{OR})} = e^{\text{LN}\left(\frac{P}{1-P}\right)} = \frac{P}{1-P}$$

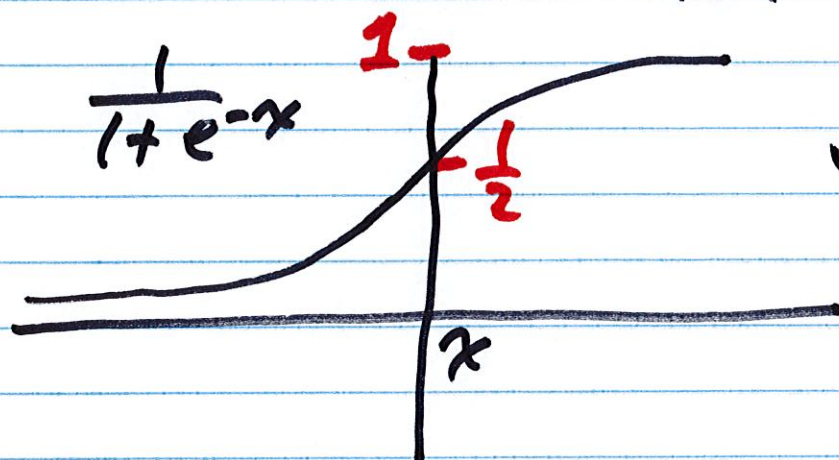
$$e^{\text{LOR}} = \frac{P}{1-P}$$

$$e^{-\text{LOR}} = \frac{1-P}{P} = \frac{1}{P} - 1$$

$$1 + e^{-\text{LOR}} = \frac{1}{P}$$

$$P = \frac{1}{1 + e^{-\underline{\text{LOR}}}}$$

$f(x) = \frac{1}{1 + e^{-x}}$ IS CALLED
THE "LOGISTIC
FUNCTION"



"S CURVE"

"MODEL" IS THAT

$$P = P(Y=1) = \frac{1}{1 + e^{-\text{LOR}}}$$

$$\text{LOR} = \beta_0 + \beta_1 X_1 + \dots + \beta_K X_K$$

GIVEN DATA, COMPUTE

$$\widehat{\text{LOR}}_i = b_0 + b_1 X_{i1} + \dots + b_K X_{iK}$$

"LOR-HAT"

WHERE X_{ij} IS VALUE OF j TH INDEX VAR FOR OBSERVATION i .

$$\hat{p}_i = \frac{1}{1 + e^{-\widehat{\text{LOR}}_i}} = \frac{\widehat{\text{OR}}_i}{1 + \widehat{\text{OR}}_i}$$

NOTE: e^{b_i} GIVES FACTOR CHANGE IN OR FOR UNIT INCREASE IN X_i

Ex: $e^{0.5} = 1.0513$ $e^{.40} = 1.492$ [$\approx +50\%$]
 $e^{.70} = 2.013$ [$\approx 100\%$]

MIN SSE DOES NOT MAKE AT
A CRITERION FOR CHOOSING
EQUATION $[b_0, b_1, \dots, b_k]$

INSTEAD, CHOOSE b_0, b_1, \dots, b_k TO
MAXIMIZE LIKELIHOOD

$$\left[\prod_{i|y_i=1} \hat{p}_i \right] \cdot \left[\prod_{i|y_i=0} (1 - \hat{p}_i) \right]$$

PRODUCT SUCH THAT

- ~~THIS~~ THIS CAN BE DONE!
BY ANY STATS PACKAGE
(NOT EXCEL)
- IN ADDITION TO POINT ESTIMATES
 b_0, b_1, \dots, b_k CAN GET CONF
INTERVALS \Rightarrow EVALUATE
SIGNIFICANCE OF VARIABLES

WHAT ABOUT RESIDUALS?

PROBLEM: GET SOMEONE LOOKING AT INDIVIDUAL OBSERVATIONS

SOLUTION: AGGREGATE DATA POINTS INTO GROUPS. CUMULATE OBSERVED # WITH $Y_i = 1$ TO EXPECTED #.

→ LETS NATURALLY TO CHI-SQUARE TESTS ON RESIDUALS. SEVERAL VERSIONS USED IN STD. SOFTWARE.

WOULD LIKE TO VISUALIZE RESIDUALS FOR A CONTINUOUS X_j (AGE, AVE-AGE)

- GROUP VARIABLE X_j INTO RANGES
- FOR EACH RANGE, CUMULATE OBS. # OF $Y_i = 1$ TO EXPECTED #

8

USING REGRESSION, CONSIDER
EACH Y_i TO BE OUTCOME
OF A BERNOULLI TRIAL
WITH $P(Y_i=1) = \hat{p}_i$
 \Rightarrow VARIANCE = $\hat{p}_i(1-\hat{p}_i)$

SO FOR A GROUP, COMPUTE
 $\sum Y_i$ WITH $\sum \hat{p}_i$
STANDARDISE USING VARIANCE...

$$\Rightarrow \text{"PEARSON RESIDUAL"} \quad \frac{\sum Y_i - \sum \hat{p}_i}{\sqrt{\sum \hat{p}_i(1-\hat{p}_i)}}$$

[SUMS ARE FOR OBSERVATIONS
IN THE GROUP]

EX: LOOK AT PEARSON RESIDUALS
AGAINST AGE AND AVE-AGE
IN CC DATA.