# Logistic Regression w/ Ungrouped Data.

In OLS, grouped vs. ungrouped don't really matter.

Ungrouped data: either having at least one continuous or interval covariate OR having a large # of categorical covariates (such that each level combination has almost no data).

Example: Age is commonly ungrouped

⤷ Artificially group data. This is often unsatisfactory.
  1) Analysis hinges dramatically on grouping.
  2) Can have really small counts.

Another example: Logistic regression is commonly used to analyze survey data. Typical survey might have 50 questions, each Q has 5 levels. Even w/ large sample, many combinations have zeros.

KEY DIFFERENCE IN ANALYZING GROUPED
vs. UNGROUPED DATA (IN LOGISTIC REG.)
IS THAT THE DEVIANCE MEANS DIFF. THINGS.
IN PARTICULAR, IN UNGROUPED DATA, THE
DEVIANCE BY ITSELF DOESN'T MEAN ANYTHING.

$\hookrightarrow$ THIS MEANS WE HAVE TO DO GOODNESS-OF-
FIT IN A DIFFERENT WAY.

ONE WAY: GROUP DATA IN SOME MANNER !
USE PEARSON $\chi^2$ OR DEVIANCE
GOODNESS-OF-FIT.

ANOTHER WAY: GROUP DATA ACCORDING TO ESTIMATED
PROBABILITIES. KNOWN AS "HOSMER-LEMESHOW"

PROCEDURE: 1) FIT SENSIBLE MODEL USING
CONTINUOUS COVARIATES

EXAMPLE: 2) FORM $K$ GROUPS OF (APPROX.) EQUAL
$K = 10$    SIZE. IF THERE $\Lambda$ OBSERVATIONS,
$\Lambda/K$ IN GROUP 1 THAT CORRESPOND TO
HIGHEST FITTED PROBABILITY. THEN
GROUP 2 WOULD HAVE NEXT $\Lambda/K$
HIGHEST, AND SO ON.

FORM $\quad H = \sum_{k=1}^{K} \frac{(O_k - E_k)^2}{\Lambda_k \bar{\pi}_k (1 - \bar{\pi}_k)}$

WHERE: $O_k$ = # OF OBSERVED COUNTS $\left( \begin{array}{c} \text{TOTAL \# OF} \\ \text{OBS IN } K^{TH} \\ \text{GROUP} \end{array} \right)$

$E_k$ = # OF EXPECTED COUNTS $\left( \frac{\Lambda}{K} + 0.5 \right)$

$\Lambda_k$ = # OBS IN GROUP $K$

$\bar{\pi}_k$ = AVERAGE OF $\hat{\pi}_i$ IN GROUP $K$.

$H \sim \chi^2_{K-2} \quad \rightarrow$ LARGE $H$ (SMALL $p$-VAL)
INDICATE MODEL MISFIT.

# Logistic Regression for Multi Level Responses

Multinomial Logistic Regression.

Example    Suppose each subject falls into one of three levels
$$\begin{cases} (1) \text{ HIV neg.} \\ (2) \text{ HIV pos. / no AIDS} \\ (3) \text{ HIV pos. / AIDS} \end{cases}$$

Let $\pi_1, \pi_2, \pi_3$ be probabilities of being on each level. We would like to know the effect of these explanatory variables on these probabilities.

$\hookrightarrow$ Note: $\pi_1 + \pi_2 + \pi_3 = 1$, so any two determine third one.

Form:
$$\theta_2 = \log\left(\frac{\pi_2}{\pi_1}\right), \quad \theta_3 = \log\left(\frac{\pi_3}{\pi_1}\right)$$

Where: $\theta_2$ is log odds of being HIV pos / no AIDS relative to HIV neg.

Idea: Fit separate $\overset{\log}{\text{odds-ratios}}$ w/ different ~~logistic~~ regressions.
        linear

Example: Suppose we have explanatory var. # of sexual partners in last 6 months (X)
$$\log\left(\frac{\pi_2(x)}{\pi_1(x)}\right) = \beta_0 + \beta_1 X$$
$$\log\left(\frac{\pi_3(x)}{\pi_1(x)}\right) = \alpha_0 + \alpha_1 X$$

$\hookrightarrow$ From here, everything is the same as before.

Caution: Be sure you know what the reference level is.

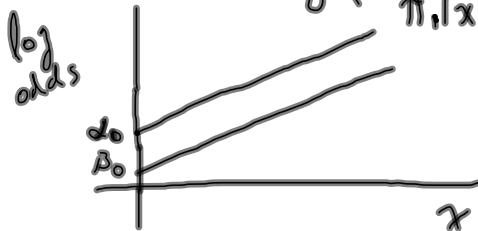## ORDERED LEVELS:

"NO" DISEASE, "SOME" DISEASE, "SEVERE" DISEASE.

THERE ARE A FEW WAYS OF GOING ABOUT THIS.

1) "COMMON SLOPES"

$$\log\left(\frac{\pi_2(x)}{\pi_1(x)}\right) = \beta_0 + \beta_1 x$$

$\uparrow$ DIFF    $\{$ SAME

$\left(\begin{array}{c}\text{WE EXPECT}\\ \text{TO SEE}\\ \hat{\beta}_1 > 0\end{array}\right)$

$$\log\left(\frac{\pi_3(x)}{\pi_1(x)}\right) = \alpha_0 + \beta_1 x$$

"log ODDS OF HIV POS/AIDS NO RELATIVE TO HIV NEG IS A CONSTANT SHIFT OF HIV POS/AIDS YES RELATIVE TO HIV NEG."



log odds

$\alpha_0$
$\beta_0$

$x$

"PROPORTIONAL ODDS MODEL"

HERE, WE LOOK AT CUMULATIVE PROBABILITY AS THE RESPONSE. DEFINE $\gamma_1 = \pi_1$, $\gamma_2 = \pi_1 + \pi_2$, $\gamma_3 = \pi_1 + \pi_2 + \pi_3 = 1$. DO INFERENCE ON $\gamma_1, \gamma_2$.

NOW, 
$$\log\left(\frac{\gamma_{(x)}}{1-\gamma_{(x)}}\right) = \beta_0 + \beta_1 x$$

$\uparrow$ DIFF    $\{$ SAME

$\left(\begin{array}{c}\text{EXPECT}\\ \hat{\beta}_1 < 0\end{array}\right)$

$$\log\left(\frac{\gamma_2(x)}{1-\gamma_2(x)}\right) = \alpha_0 + \beta_1 x$$

WE CAN USE OTHER LINK FUNCTIONS
THAT LOOK/ACT A LOT LIKE THE LOGIT.

ANY $F(\pi)$ THAT IS
    1) DIFFERENTIABLE
    2) GOES FROM $-\infty, \infty$ AS $\pi$ GOES $0 \to 1$
    3) ALWAYS INCREASING
CAN BE USED AS A LINK.
    EXAMPLES
- LOGISTIC REGRESSION W/ CANONICAL LINK.

- USE THE INVERSE CDF OF A STANDARD NORMAL
  AS THE LINK. GIVES "PROBIT" MODEL.