

Unit 7

Logistic Regression

“To all the ladies present and some of those absent”

- Jerzy Neyman

What behaviors influence the chances of developing a sexually transmitted disease? Comparing demographics, health education, access to health care, which of these variables are significantly associated with failure to obtain an HIV test? Among the several indicators of risk, including age, co-morbidities, severity of disease, which are significantly associated with surgical mortality among patients undergoing transplant surgery?

In all of these examples, the outcome observed for each individual can take on only one of two possible values: positive or negative test, alive or dead, remission or non-remission, and so on. Collectively, the data to be analyzed are **proportions**.

Proportions have some important features that distinguish them from data measured on a continuum. Proportions (1) are **bounded** from below by the value of zero (or zero percent) and bounded from above by one (or 100 percent); (2) as the proportion gets close to either boundary, the variance of the proportion gets smaller and smaller; thus, we **cannot assume a constant variance**; and (3) proportions are **not distributed normal**. **Normal theory regression models are not appropriate for the analysis of proportions.**

In unit 4, Categorical Data Analysis, emphasis was placed on contingency table approaches for the analysis of such data. It was highlighted that these methods should always be performed for at least two reasons: (1) they give a good feel for the data; and (2) they are free of the assumptions required for regression modeling.

Unit 7 is an introduction to logistic regression approaches for the analysis of proportions where it is of interest to explore the roles of possibly several influences on the observed proportions.

Table of Contents

Topic		
	Learning Objectives	3
	1. From Linear Regression to Logistic Regression	4
	2. Use of VDT's and Spontaneous Abortion	5
	3. Definition of the Logistic Regression Model	7
	4. Estimating Odds Ratios	11
	5. Estimating Probabilities	16
	6. The Deviance Statistic	17
	a. The Likelihood Ratio Test	20
	b. Model Development	23
	7. Illustration – Depression Among Free-Living Adults	26
	8. Regression Diagnostics	36
	a. Assessment of Linearity	39
	b. Hosmer-Lemeshow Goodness of Fit Test	40
	c. The Linktest	43
	d. The Classification Table	45
	e. The ROC Curve	47
	f. Pregibon Delta Beta Statistic	49
	9. Example - Disabling Knee Injuries in the US Army	50
Appendix	Overview of Maximum Likelihood Estimation	58

Learning Objectives

When you have finished this unit, you should be able to:

- Explain why a normal theory regression model is *not* appropriate for a regression analysis of proportions.
- State the *expected value* (the mean) of a Bernoulli random variable.
- Define the *logit* of the mean of a Bernoulli random variable.
- State the *logistic regression model* and, specifically, the logit link that relates the logit of the mean of a Bernoulli random variable to a linear model in the predictors.
- Explain how to *estimate odds ratio measures* of association from a fitted logistic regression model.
- Explain how to *estimate probabilities of event* from a fitted logistic regression model.
- Perform and interpret *likelihood ratio test* comparisons of hierarchical models.
- Explain and compare *crude versus adjusted* estimates of odds ratio measures of association.
- Assess *confounding* in logistic regression model analyses.
- Assess *effect modification* in logistic regression model analyses.
- *Draft an analysis plan* for multiple predictor logistic regression analyses of proportions.

1. From Linear Regression To Logistic Regression An Organizational Framework

In **unit 5** (*Regression and Correlation*), we considered single and multiple predictor regression models for a single outcome random variable Y assumed **continuous** and distributed **normal**.

In **unit 7** (*Logistic regression*), we consider single and multiple regression models for a single outcome random variable Y assumed discrete, **binary**, and distributed **bernoulli**.

	Unit 5 Normal Theory Regression	Unit 7 Logistic Regression
Y	<ul style="list-style-type: none"> - univariate - continuous - Example: $Y = \text{cholesterol}$ 	<ul style="list-style-type: none"> - univariate - discrete, binary - Example: $Y = \text{dead/alive}$
X_1, X_2, \dots, X_p	<ul style="list-style-type: none"> - one or multiple - discrete or continuous - treated as fixed 	<ul style="list-style-type: none"> - one or multiple - discrete or continuous - treated as fixed
$Y \mid X_1=x_1, \dots, X_p=x_p$	- Normal (Gaussian)	- Bernoulli (or binomial)
$E(Y \mid X_1=x_1, \dots, X_p=x_p)$	$\mu_{Y X_1 \dots X_p} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$	$\mu_{Y X_1 \dots X_p} = \pi_{Y X_1 \dots X_p}$ $= \frac{1}{1 + \exp[-(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)]}$
Right hand side of model	$\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$	$\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$
Link	"natural" or "identity" $\mu_{Y X_1 \dots X_p}$ $= \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$	"logit" $\text{logit}[\mu_{Y X_1 \dots X_p}]$ $= \text{logit}[\pi_{Y X_1 \dots X_p}]$ $= \ln \left[\frac{\pi_{Y X_1 \dots X_p}}{1 - \pi_{Y X_1 \dots X_p}} \right]$ $= \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$
Estimation	Least squares (= maximum likelihood)	Maximum Likelihood
Tool	Residual sum of squares	Deviance statistic
Tool	Partial F Test	Likelihood Ratio Test

2. Use of Video Display Terminals and Spontaneous Abortion

Data Used.

vdt.Rdata

Example.

Source: Schnorr et al (1991) Video Display Terminals and the Risk of Spontaneous Abortion. *New England Journal of Medicine* 324: 727-33.

Background:

Adverse pregnancy outcomes were correlated with use of video display terminals (VDT's) beginning in 1980.

Subsequent studies were inconsistent in their findings.

Previous exposure assessments were self-report or derived from job title descriptions.

Electromagnetic fields were not previously measured.

Research Question:

What is the nature and significance of the association, as measured by the **odds ratio**, between exposure to electromagnetic fields emitted by VDTs and occurrence of spontaneous abortion, after controlling for

- History of prior spontaneous abortion
- Cigarette Smoking
- History of thyroid condition

Design: Retrospective cohort investigation of two groups of full-time female telephone operators.

882 Pregnancies:	N	<u>Spontaneous Abortion</u>	
		n	%
Exposed	366	54	14.8%
Unexposed	516	82	15.9%

The Data (vdt.Rdata):

Variable	Label	Range/Codes
AVGVDT	average hours vdt in 1st trimester	continuous
NUMCIGS	# cigarettes/day	continuous
PRIORSAB	prior spontaneous abortion	1=yes, 0=no
SAB	spontaneous abortion	1=yes, 0=no
SMOKSTAT	smoker	1=yes, 0=no
PRTHYR	prior thyroid condition	1=yes, 0=no
VDTEXPOS	VDT exposure	1=yes, 0=no

<u>AVGVDT</u>	<u>NUMCIGS</u>	<u>PRIORSAB</u>	<u>SAB</u>	<u>SMOKSTAT</u>	<u>PRTHYR</u>	<u>VDTEXPOS</u>
0.000	15	0	0	1	0	0
0.000	10	0	0	1	0	0
0.000	20	0	0	1	0	0
	20	0	0	1	0	1
27.764	20	0	1	1	0	1
28.610	0	0	0	0	0	1
0.000	0	0	0	0	0	0
	0	0	0	0	0	1
19.717	0	0	0	0	0	1
0.000	0	0	0	0	0	0
25.022	0	0	0	0	0	1
...
0.000	0	1	0	0	0	0

3. Definition of the Logistic Regression Model

Recall. We suspect that multiple factors, especially use of video display terminals, contribute to an individual's odds of event of spontaneous abortion.

The outcome or dependent event variable is **Y=sab**. Its value is y and

$$\begin{aligned} &= 1 \text{ if spontaneous abortion occurred} \\ &0 \text{ otherwise} \end{aligned}$$

The predictors that might influence the odds of SAB are several:

$$\begin{aligned} X_1 &= \text{avgvdt} \\ X_2 &= \text{numcigs} \\ X_3 &= \text{priorsab} \\ X_4 &= \text{smokstat} \\ X_5 &= \text{prthyr, and} \\ X_6 &= \text{vdtexpos} \end{aligned}$$

We are especially interested in exposure for which we have two predictors, X_6 and X_1 :

$$\begin{aligned} X_6 &= \text{vdtexpos (coded = 1 for exposed and = 0 for NON exposed) and} \\ &= 1 \text{ if exposed} \\ &0 \text{ for NON exposed; and} \\ X_1 &= \text{avgvdt} \end{aligned}$$

Among the $N=882$ in our sample, we have potentially $N=882$ unique probabilities of spontaneous abortion.

$$\pi_1, \pi_2, \dots, \pi_N.$$

For the i^{th} person

$$\pi_i = \text{Function (} X_{1i}, X_{2i}, X_{3i}, X_{4i}, X_{5i}, X_{6i} \text{) defined}$$

$$\begin{aligned} \pi_i &= \Pr [Y_i = 1] \\ (1 - \pi_i) &= \Pr [Y_i = 0] \end{aligned}$$

So, how do we model the $N=882$ individual probabilities π_i in relationship to the predictors?

Recall. Each profile of values, $\underline{X} = [X_1=x_1, X_2=x_2, \dots, X_6=x_6]$, defines a sub-population with their own distribution of outcomes Y . For example the women with $X_3=1$ are the women with a history of prior spontaneous abortion, and are distinct from the women with $X_3=0$ (who have no such prior history). And so on; we can talk about distinct sub-populations based on the entire profile of values on X_1, X_2, \dots, X_6 .

Review of [normal theory](#) linear regression analysis:

$Y \mid [X_1, X_2, X_3, X_4, X_5, X_6]$ (read: “ Y given $[X_1, X_2, X_3, X_4, X_5, X_6]$ ” is assumed to be distributed normal (Gaussian)

with mean $= \mu_{Y|\underline{x}}$ and variance $= \sigma_{Y|\underline{x}}^2$.

The mean of Y at $[X_1, X_2, X_3, X_4, X_5, X_6]$ is modeled linearly in $\underline{x} = [X_1, X_2, X_3, X_4, X_5, X_6]$

Thus, the mean of $Y \mid [X_1, X_2, X_3, X_4, X_5, X_6]$ $= E[Y \mid (X_1, X_2, X_3, X_4, X_5, X_6)] = \mu_{Y|\underline{x}}$

In normal theory linear regression:

$$E[Y \mid \underline{x}] = \mu_{\underline{x}} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6$$

↑
“natural link”

↑
“right hand side is linear in the predictors”

In a [logistic](#) model regression analysis, the framework is a little different:

Y is assumed to be distributed Bernoulli

with mean $= \pi_{\underline{x}}$ and variance $= \pi_{\underline{x}}(1-\pi_{\underline{x}})$

We do not model the mean of $Y \mid \underline{X}=\underline{x} = \pi_{\underline{x}}$ linearly in $\underline{x} = [X_1 \dots X_6]$.

Instead, we model the logit of the mean of $Y \mid \underline{X}=\underline{x} = \pi_{\underline{x}}$ linearly in $\underline{x} = [X_1 \dots X_6]$.

$$\text{Logit} [E(Y \mid \underline{X})] = \text{logit} [\pi_{\underline{x}}] = \ln \left[\frac{\pi_{\underline{x}}}{1-\pi_{\underline{x}}} \right] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_5 X_5 + \beta_6 X_6$$

↑
“logit link of expected Y ”

↑
“right hand side is linear in the predictors”

Solution for Probability $[Y=1 | X_1=x_1, X_2=x_2, \dots, X_6=x_6] = E[Y | X_1=x_1, X_2=x_2, \dots, X_6=x_6]$:

The formula for $\Pr [Y = 1 | X_1=x_1, X_2=x_2, \dots, X_6=x_6]$ can be written in either of two ways:

$$\begin{aligned}\pi_x &= \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_6 x_6)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_6 x_6)} \\ &= \frac{1}{1 + \exp[-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_6 x_6)]}\end{aligned}$$

$\Pr [Y = 0 | X_1=x_1, X_2=x_2, \dots, X_6=x_6]$ is

$$(1 - \pi_x) = \frac{1}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_6 x_6)}$$

The logistic regression model has other names, too (sorry!): “log-linear odds” and “exponential odds”

The reason this makes any sense at all is the following. Logistic regression model focuses on the odds of event (in this case event of spontaneous abortion, SAB).

$$1) \quad \ln [\text{odds} (\pi_x)] = \ln \left[\frac{\pi_x}{1 - \pi_x} \right] = \beta_0 + \dots + \beta_6 x_6 \text{ is a } \text{log-linear odds} \text{ model.}$$

$$2) \quad \left[\frac{\pi_x}{1 - \pi_x} \right] = \exp \{ \beta_0 + \dots + \beta_6 x_6 \} \text{ is an } \text{exponential odds} \text{ model.}$$

We do *not* model $E[Y | X] = \pi_x = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6$!!

$$1) \quad \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6$$

can range from $-\infty$ to $+\infty$ but π_x ranges from 0 to 1.

$$2) \quad \pi_x = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6 \text{ is often not a good description of nature.}$$

Assumptions in a logistic regression model:

- 1) Each Y_i follows a distribution that is **Bernoulli** with parameter $E[Y | X] = \pi_{x_i}$.
- 2) The Y_1, Y_2, \dots, Y_N are independent.
- 3) The values of the predictors, $X_{i1}=x_{i1} \dots X_{i6}=x_{i6}$, are treated as fixed.
- 4) The model is correct (this is also referred to as “**linearity in the logit**”).

$$\begin{aligned}
 \text{logit}[E(Y) | X_1=x_1, X_2=x_2, \dots, X_6=x_6] \\
 &= \text{logit}[\pi_x] \\
 &= \ln \left[\frac{\pi_x}{1 - \pi_x} \right] \\
 &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \beta_6 X_6
 \end{aligned}$$

- 5) No multicollinearity
- 6) No outliers
- 7) Independence

3. Estimating Odds Ratios

The goal is to understand a fitted model. We'll get to the details of estimation later. Once a logistic regression model has been fit, the prediction equation can be used to estimate **odds ratio (OR) measures of association**.

Example 1: What is the estimated crude relative odds (OR) of spontaneous abortion (SAB) associated with any exposure (1 = exposed, 0 = not exposed) to a video display terminal (VDTEXPOS)?

Step 1:

To obtain crude odds ratios, either a 2x2 table can be used or a one predictor logistic regression model can be fit. Here, it is given by

$$\text{logit} \{ \text{probability} [\text{SAB}=1] \} = \beta_0 + \beta_1 \text{VDTEXPOS for}$$

$$\text{VDTEXPOS} = 1 \text{ if exposed} \\ 0 \text{ if NOT exposed}$$

R

```
setwd("/Users/cbigelow/Desktop/")
options(scipen=1000)
options(signif.stars=FALSE)
load(file="vdt.Rdata")

# display 2x2 table
library(summarytools)
library(stargazer)
summarytools::cTable(vdtdata$vdtxpos, vdtdata$sab, prop = 'n', totals = TRUE)

## Cross-Tabulation
## Variables: vdtxpos * sab
## Data Frame: vdtdata
##
## -----
##      sab      0      1      Total
##  vdtxpos
##      0      434      82      516
##      1      312      54      366
##      Total      746     136     882
## -----

# Fit model
fit1 <- glm(sab ~ vdtxpos, data=vdtdata, family=binomial)
summary(fit1)
## Call:
## glm(formula = sab ~ vdtxpos, family = binomial, data = vdtdata)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.5883  -0.5883  -0.5650  -0.5650   1.9564
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -1.66633     0.12041 -13.838 <0.000000000000002 ***
## vdtxpos       -0.08769     0.19032  -0.461      0.645
```

Yielding the following prediction equation

$$\text{fitted logit} \{ \text{probability}[\text{sab}=1] \} = -1.66633 - 0.08769 \cdot \text{vdtxpos}$$

Nature ——— Population/ ——— Observation/ ——— Relationships/ ——— Analysis/
Sample Data Modeling Synthesis

Step 2:

Recognize a wonderful bit of algebra.

For a single exposure variable (1=exposed, 0=not exposed)

$$\begin{aligned} \text{OR}_{1 \text{ versus } 0} &= \exp\{\beta\} \text{ where } \beta = \text{regression parameter for the exposure variable} \\ &= \exp\{\text{logit}(\pi_1) - \text{logit}(\pi_0)\} \end{aligned}$$

Proof (read if you are interested!):

$$\begin{aligned} \text{OR} &= \exp\{\ln[\text{OR}]\} \\ &= \exp\left\{\ln\left[\frac{\pi_1/(1-\pi_1)}{\pi_0/(1-\pi_0)}\right]\right\} \\ &= \exp\left\{\ln\left[\frac{\pi_1}{1-\pi_1}\right] - \ln\left[\frac{\pi_0}{1-\pi_0}\right]\right\} \\ &= \exp\{\text{logit}(\pi_1) - \text{logit}(\pi_0)\} \end{aligned}$$

“1” is the comparison and is vdtexpos=1:

$$\begin{aligned} \text{Estimated logit}\{\text{prob}[SAB=1 \mid \text{vdtexpos}=1]\} &= \hat{\beta}_0 + \hat{\beta}_1 \\ &= -1.66633 - 0.08769 \end{aligned}$$

“0” is the reference and is vdtexpos=0:

$$\text{Estimated logit}\{\text{prob}[SAB=1 \mid \text{vdtexpos}=0]\} = \hat{\beta}_0 = -1.66633$$

Step 3: Apply.

The odds ratio measure of association comparing the exposed telephone operator (“1”) to the unexposed telephone operator (“0”) is

$$\begin{aligned} \text{OR} &= \exp\{\text{logit}(\pi_1) - \text{logit}(\pi_0)\} \\ &= \exp\{[\beta_0 + \beta_1] - [\beta_0]\} \\ &= \exp\{\beta_1\} \\ &= \exp\{-0.08769\} \\ &= 0.9160 \rightarrow \text{“Compared to the unexposed, the exposed have a relative odds of spontaneous abortion=.916”} \end{aligned}$$

The two profiles being compared can differ on several predictors! Let's try another one.

Here is the wonderful algebra, in all its glory:

For two profiles of predictor variable values, “comparison” versus “reference”

$$\text{OR}_{\text{comparison versus reference}} = \exp \{ \text{logit} (\pi_{\text{comparison}}) - \text{logit} (\pi_{\text{reference}}) \}$$

$$\ln \{ \text{OR}_{\text{comparison versus reference}} \} = \text{logit} (\pi_{\text{comparison}}) - \text{logit} (\pi_{\text{reference}})$$

Example 2 - For the fitted model below, calculate the estimated relative odds (OR) of spontaneous abortion (SAB) for a person who is not exposed to a VDT, smokes 10 cigarettes per day, has no history of prior SAB, and no thyroid condition relative to a person who has an average of 20 hours exposure to a VDT, is a nonsmoker, has a history of prior SAB and does have a thyroid condition?

$$\begin{aligned} \text{fitted logit} \{ \text{prob}[sab=1] \} \\ = -1.95958 + 0.00508(\text{avgvdt}) + 0.04267(\text{numcigs}) + 0.38500(\text{priorsab}) + 1.27420(\text{prthyrr}) \end{aligned}$$

Step 1:

Calculate the two predicted logits and compute their difference.

	Value of Predictor for Person	
	“comparison”	“reference”
avgvdt	0	20
numcigs	10	0
priorsab	0	1
prthyrr	0	1

“comparison”

$$\begin{aligned} \text{fitted logit} [\pi_{\text{comparison}}] &= -1.95958 + 0.00508(0) + 0.04267(10) + 0.38500(0) + 1.27420(0) \\ &= -1.5329 \end{aligned}$$

“reference”:

$$\begin{aligned} \text{fitted logit} [\pi_{\text{reference}}] &= -1.95958 + 0.00508(20) + 0.04267(0) + 0.38500(1) + 1.27420(1) \\ &= -0.1988 \end{aligned}$$

$$\begin{aligned} \text{logit} [\pi_{\text{comparison}}] - \text{logit} [\pi_{\text{reference}}] &= -1.5329 - [-0.1988] \\ &= -1.3341 \end{aligned}$$

Step 2:

Exponentiate.

$$\begin{aligned} \text{OR}_{\text{comparison versus reference}} &= \exp \{ \text{logit} [\pi_{\text{comparison}}] - \text{logit} [\pi_{\text{reference}}] \} = \exp \{ -1.3341 \} \\ &= \mathbf{0.2634} \end{aligned}$$

Interpretation - Comparing two odds: The estimated odds of spontaneous abortion (SAB) for a person who is not exposed to a VDT, smokes 10 cigarettes per day, has no history of prior SAB, and no thyroid condition is 0.2631 times that of the odds of spontaneous abortion (SAB) for a person who has an average of 20 hours exposure to a VDT, is a nonsmoker, has a history of prior SAB and does have a thyroid condition. In odds ratio (OR) parlance: The relative odds (odds ratio OR) is .2631

R Illustration

```
# Fit model
fit2 <- glm(sab ~ avgvdt + numcigs + priorsab + prthyr,
            family=binomial,
            data=vdtdata)

# Define profiles to be compared
comparison <- data.frame(avgvdt=0,numcigs=10,priorsab=0,prthyr=0)
referent <- data.frame(avgvdt=20,numcigs=0,priorsab=1,prthyr=1)

# Obtain predicted logits for each profile
logitc <- predict(fit2,comparison,type="link")
logitr <- predict(fit2,referent,type="link")

# Ln [ OR ] = difference in predicted logits
# [OR] = exp [difference in predicted logits]
exp(logitc - logitr)
## 1
## 0.263115
```

In General:

The Odds Ratio estimate (\hat{OR}) of association with outcome accompanying a unit change in the predictor X is a function of the estimated regression parameter $\hat{\beta}$

$$\hat{OR}_{\text{UNIT change in X}} = \exp \{ \hat{\beta} \}$$

Tip – $OR_{10 \text{ unit change in X}} = \exp [10 * \beta]$

The hypothesis test of null: $OR=1$
is equivalent to

The hypothesis test of null: $\beta = 0$

For a rare outcome (typically disease), the relative risk (\hat{RR}) estimate of association with outcome accompanying a unit change in the predictor X is reasonably estimated as a function of the estimated regression parameter β

$$\hat{RR}_{\text{UNIT change in X}} = \exp \{ \hat{\beta} \}, \text{ approximately}$$

5. Estimating Probabilities

The goal is to understand a fitted model.

Once a logistic regression model has been fit, the prediction equation can also be used to estimate **probabilities of event occurrence**. The prediction equation can be used to estimate probabilities of event of **disease** if the study design is a **cohort**; it is used to estimate probabilities of **history of exposure** if the study design is **case-control**.

*Reminder ... it is **not possible to estimate probability of disease from analyses of case-control studies**.*

Recall that for Y distributed Bernoulli, its expected value is given by

$$E[Y] = \pi = \text{Probability of event occurrence}$$

Example 1- Consider again the one predictor model below. Under the assumption of a cohort study design, what is estimated probability of spontaneous abortion (sab) for a person with any exposure to a video display terminal? Consider the single predictor model containing the predictor vdtexpos)

$$\text{Predicted logit } \{ \text{prob}[SAB=1 | \text{vdtexpos}] \} = \hat{\beta}_0 + \hat{\beta}_1[\text{vdtexpos}] = -1.66633 - 0.08769 \cdot \text{VDTEXPOS}$$

Step 1:

Utilizing the algebra on page 9, we have:

$$\text{Estimated pr}[SAB=1] = \hat{\pi}_{\text{VDTEXPOS}=1} = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1[\text{vdtexpos}])}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1[\text{vdtexpos}])} = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1)}$$

Step 2:

Set VDTEXPOS=1, $\beta_0 = -1.66633$, $\beta_1 = -0.08769$ and solve

$$\begin{aligned} \text{Estimated pr}[SAB=1] &= \frac{\exp(-1.66633 - 0.08769[1])}{1 + \exp(-1.66633 - 0.08769[1])} \\ &= \frac{0.1731}{1.1731} = 0.148 \end{aligned}$$

R

```
fit1 <- glm(sab ~ vdtexpos, data=vdtdata, family=binomial)
x_exposed <- data.frame(vdtexpos=1)
phat_exposed <- round(predict(fit1,x_exposed,type="response"), 3)

paste("Predicted Probability of SAB for (VDTEXPOS=1) = ",phat_exposed)
## [1] "Predicted Probability of SAB for (VDTEXPOS=1) = 0.148"
```

6. The Deviance Statistic

"G Statistic", "Log likelihood Statistic", "Scaled Deviance", Residual Deviance"

The goal is to understand deviance (hint: think sums of squares in an analysis of variance). Recall the concept of "analysis of variance" introduced in Unit 5, Regression and Correlation. Analysis of variance is about the total variability of the observed outcome, and its partitioning into portions that are explained by the fitted model (due model/due regression) versus what's left over as unexplained (due residual/due error). The deviance statistic in logistic regression is a measure of what remains left over as unexplained by the fitted model, analogous to the residual sum of squares in normal theory regression.

But first, a few words about likelihood, L .

$L_{\text{saturated}}$: We get the largest likelihood of the data when we fit a model that allows a separate predictor for every person. This is called the likelihood of the saturated model.

$L_{\text{saturated}}$ is a large number.

L_{current} : We get an estimated likelihood of the data when we fit the current model.

L_{current} is a smaller number.

The deviance statistic in logistic regression is related to the two likelihoods, L_{current} and $L_{\text{saturated}}$ in the following way.

The current model explains a lot	The current model does NOT explain a lot
$L_{\text{current}} \approx L_{\text{saturated}}$	$L_{\text{current}} \ll L_{\text{saturated}}$
$\frac{L_{\text{current}}}{L_{\text{saturated}}} \approx 1$	$\frac{L_{\text{current}}}{L_{\text{saturated}}} \ll 1$
$\ln \left[\frac{L_{\text{current}}}{L_{\text{saturated}}} \right] \approx 0$	$\ln \left[\frac{L_{\text{current}}}{L_{\text{saturated}}} \right] \ll 0$
Deviance = $(-2) \ln \left[\frac{L_{\text{current}}}{L_{\text{saturated}}} \right] \approx 0$	Deviance = $(-2) \ln \left[\frac{L_{\text{current}}}{L_{\text{saturated}}} \right] \gg 0$
A number close to 0	A large positive number

Evidence that the current model explains a lot of the variability in outcome

Deviance \approx small

p-value \approx large

$$\begin{aligned}\text{Deviance Statistic, } D &= -2 \ln \left[\frac{L_{\text{current}}}{L_{\text{saturated}}} \right] \\ &= (-2) \ln (L_{\text{current}}) - (-2) \ln (L_{\text{saturated}})\end{aligned}$$

$$\text{Deviance df} = [\text{Sample size}] - [\# \text{ fitted parameters}]$$

where

$$\begin{aligned}L_{\text{current}} &= \text{likelihood of data using current model} \\ L_{\text{saturated}} &= \text{likelihood of data using the saturated model}\end{aligned}$$

Notes -

- (1) By itself, the deviance statistic does not have a well-defined distribution
- (2) However, **differences of deviance statistics** that compare hierarchical models do have well defined distributions; when the null hypothesis is true, they are distributed chi square.

A Feel for the Deviance Statistic

- (1) Roughly, the **deviance statistic D** is a measure of what remains unexplained.
Hint – The analogue in normal theory regression is the residual sum of squares (SSQ error)
- (2) A deviance statistic value **close to zero** says that a lot is explained and, importantly, that little remains unexplained. → The current model with its few predictors performs similarly to the saturated model that permits a separate predictor for each person.
- (3) **WARNING!** The deviance statistic D is **NOT** a measure of goodness-of-fit. Recall that we said the same thing about the overall F-statistic in normal theory regression.
- (4) The **deviance statistic D** is the basis of the likelihood ratio test.
- (5) The likelihood ratio test is used for the comparison of hierarchical models.
Recall – In normal theory regression, hierarchical models are compared using the Partial F-test.

R

In R, “*deviance residual*” = $(-2) \times \ln\text{-Likelihood}$.

```
# Fit model
fit1 <- glm(sab ~ vdtexpos,
            data=vdtdata,
            family=binomial)

# Obtain value of ln-likelihood: logLik()
ln_likelihood <- round(logLik(fit1),5)
paste("Log-Likelihood = ", ln_likelihood )
## [1] "Log-Likelihood = -379.08045"

# In R, “deviance residual” = (-2)*ln-Likelihood: Brute force
deviance_residual <- round((-2)*ln_likelihood,4)
paste("Deviance Residual in R = (-2)*ln-likelihood = ", deviance_residual)
## [1] "Deviance Residual in R = (-2)*ln-likelihood = 758.1609"
```

a. The Likelihood Ratio (LR) Test

Likelihood Ratio (LR) Test

Under the assumptions of a logistic regression model and the comparison of the hierarchical models:

$$\text{Reduced: } \text{logit}[\pi | X_1, X_2, \dots, X_p] = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$$\text{Full: } \text{logit}[\pi | X_1, X_2, \dots, X_p, X_{p+1}, X_{p+2}, \dots, X_{p+k}] = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \beta_{p+1} X_{p+1} + \dots + \beta_{p+k} X_{p+k}$$

For testing:

$$H_O: \beta_{p+1} = \beta_{p+2} = \dots = \beta_{p+k} = 0$$

$$H_A: \text{not}$$

The **Likelihood Ratio Test** Statistic LR for testing this null hypothesis is defined

$$\begin{aligned} \text{LR} &= \text{Deviance}_{\text{REDUCED}} - \text{Deviance}_{\text{FULL}} \\ &= [(-2) \ln(L)_{\text{REDUCED}} - (-2) \ln(L)_{\text{SATURATED}}] - [(-2) \ln(L)_{\text{FULL}} - (-2) \ln(L)_{\text{SATURATED}}] \\ &= [(-2) \ln(L)_{\text{REDUCED}}] - [(-2) \ln(L)_{\text{FULL}}] \end{aligned}$$

When the **null hypothesis** is true, the LR test statistic is distributed **Chi Square_{DF=k}** with **degrees of freedom = k**

Rejection of the null hypothesis occurs for

Test statistic values, LR = large
and accompanying p-value = small

Tip – In practice, we obtain LR using the 2nd formula (see hack below); it says: $\text{LR} = [(-2) \ln(L)_{\text{REDUCED}}] - [(-2) \ln(L)_{\text{FULL}}]$

Hack!!

$$\text{Deviance}_{\text{REDUCED}} - \text{Deviance}_{\text{FULL}} = [(-2) \ln(L)_{\text{REDUCED}}] - [(-2) \ln(L)_{\text{FULL}}]$$

Example: Consider the following hierarchical fitted models, full versus reduced. Controlling for prior spontaneous abortion (PRIORSAB), is 0/1 exposure to VDT statistically significantly associated with spontaneous abortion?

Reduced Model ("reference")	Full Model ("comparison")
$\text{logit } \{\text{pr } [\text{sab}=1]\} = \beta_0 + \beta_1 \text{ PRIORSAB}$	$\text{logit } \{\text{pr } [\text{sab}=1]\} = \beta_0 + \beta_1 \text{ PRIORSAB} + \beta_2 \text{ VDTEXPOS}$
$(-2) \ln L_{\text{reduced}} = 754.56$ Deviance $DF_{\text{reduced}} = 881 - 1 = 880$	$(-2) \ln L_{\text{full}} = 754.26$ Deviance $DF_{\text{full}} = 881 - 2 = 879$

H_0 : VDTEXPOS, controlling for PRIORSAB, is **not** associated with SAB
 $\beta_{\text{VDTEXPOS}} = 0$ in the model that also contains PRIORSAB

H_A : VDTEXPOS, controlling for PRIORSAB, is associated with SAB
 $\beta_{\text{VDTEXPOS}} \neq 0$ in the model that also contains PRIORSAB

In logistic regression, the likelihood ratio test comparison of two hierarchical models is similar to the partial F test comparison of hierarchical models in normal theory linear regression. We compare two models, one of which is an enhancement of the other. Recall the meaning of hierarchical models: the predictors in the reduced model are the control variables and are a subset of the predictors in the full model. Thus, the test is assessing the statistical significance of the "extra" variables in the full model, after adjustment for the control variables in the reduced model.

Step 1: Compute the change in deviance and the change in deviance df, remembering that in logistic regression the subtraction is of the form "reduced" - "full".

$$\begin{aligned}
 \text{Likelihood Ratio Test LR} &= (-2) \ln L_{\text{reduced}} - (-2) \ln L_{\text{full}} \\
 &= 754.56 - 754.26 \\
 &= 0.30
 \end{aligned}$$

$$\begin{aligned}
 \Delta \text{ Deviance Df} &= \text{Deviance } DF_{\text{reduced}} - \text{Deviance } DF_{\text{full}} \\
 &= 880 - 879 \\
 &= 1
 \end{aligned}$$

Step 2: Compute the p-value.

p-value = $\Pr [\text{Chi Square (df=1)} > 0.30] = .58$

R

```
pchisq(q=0.30, df=1, lower.tail=F)
[1] 0.583882
```

Step 3: Interpret.

The null hypothesis is not rejected. In this sample, controlling for history of prior spontaneous abortion (PRIORSAB), exposure to video display terminal (VDTEXPOS) is not statistically significantly associated with event of spontaneous abortion (p-value = .58).

R Illustration

```
library(lmtest)
reduced <- glm(sab ~ priorsab,
               data=vdtdata,
               family=binomial)

full <- glm(sab ~ priorsab + vdtexpos,
            data=vdtdata,
            family=binomial)
lmtest::lrtest(reduced, full)

## Likelihood ratio test
##
## Model 1: sab ~ priorsab
## Model 2: sab ~ priorsab + vdtexpos
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    2 -377.28
## 2    3 -377.13 1  0.2977    0.5853
```

b. Model Development

Recall from Unit 5, Regression and Correlation with apologies, the following is a duplication

There are *no* rules *nor a single best strategy*. Different study designs and research questions call for different approaches. **Tip** – Before you begin model development, make a list of your study design, research aims, outcome variable, primary predictor variables, and covariates.

As a general suggestion, the following approach has the advantages of providing a reasonably thorough **exploration of the data and relatively little risk of missing something important**.

Preliminary – Be sure you have: (1) checked, cleaned and described your data, (2) screened the data for multivariate associations, and (3) thoroughly explored the bivariate relationships.

Step 1 – Fit the “maximal” model.

The maximal model is the large model that contains all the explanatory variables of interest as predictors. This model also contains all the covariates that might be of interest. It also contains all the interactions that might be of interest. Note the amount of variation explained.

Step 2 – Begin simplifying the model.

Inspect each of the terms in the “maximal” model with the goal of removing the predictor that is the least significant. Drop from the model the predictors that are the least significant, beginning with the higher order interactions (**Tip** -interactions are complicated and we are aiming for a simple model). Fit the reduced model. Compare the amount of variation explained by the reduced model with the amount of variation explained by the “maximal” model.

If the deletion of a predictor has little effect on the variation explained
Then leave that predictor out of the model.

And inspect each of the terms in the model again.

If the deletion of a predictor has a significant effect on the variation explained ...
Then put that predictor back into the model.

Step 3 – Keep simplifying the model.

Repeat step 2, over and over, until the model remaining contains nothing but significant predictor variables.

Beware of some important caveats

- Sometimes, you will want to keep a predictor in the model regardless of its statistical significance (an example is randomization assignment in a clinical trial)
- The order in which you delete terms from the model matters!
- You still need to be flexible to considerations of biology and what makes sense.

So, compared to normal theory regression, what's new here?

In logistic regression, instead of doing a Partial F-Test, we do a **Likelihood Ratio Test**.

If the likelihood ratio statistic is statistically significant (small p-value), we say that the added variables are statistically significant after adjustment for the control variables.

Example – Depression Among Free-Living Adults.

Among free-living adults of Los Angeles County, what is the prevalence of depression and what are its correlates? In particular, in a given data set containing information on several candidate predictors, which predictors are the significant ones?

A reasonable analysis approach for *this particular example* is the following:

Step 1. Fit single predictor models. Retain for further consideration:

- Predictors with crude significance levels of association $p < .25$
- Predictors of *a priori* interest

Step 2. Evaluate candidate predictors for evidence of multicollinearity:

Step 3. Fit a multivariable model containing the “candidates” from step 1. Retain for further consideration

- Predictors with adjusted significance levels $p < .10$

Step 4. Fit the multivariable model containing the reduced set of “candidates” from step 3.

- Compare the step 3 and step 4 models using the likelihood ratio (LR) test.

Step 5. Investigate confounding. For each confounder

- Begin with the step 4 model. --- **reduced model** ---
- Fit an enhanced model that includes the suspected confounder.
Note the estimated β 's and deviance statistic values. -- **full model** --
- Assess the adjusted statistical significance of the suspected confounder using a likelihood ratio (LR) test.
- Compute relative change in the estimated β 's:

$$\Delta\hat{\beta} = \left(\frac{|\hat{\beta}_{\text{without confounder}} - \hat{\beta}_{\text{with confounder}}|}{\hat{\beta}_{\text{with confounder}}} \right) \times 100$$

Criteria for Retention of Suspected Confounder

1. Likelihood ratio (LR) test of its adjusted association is significant; **and**
2. $\Delta\beta \geq 15\%$ or so.

Step 6. Investigate effect modification

- Begin with the “near final” model identified in step 5
- Fit, one at a time, enhanced models that contain each pairwise interaction
- Assess statistical significance of each interaction using the LR test

7. Illustration

Depression Among Free-Living Adults

Data

depress_small.Rdata

depress.Rdata

Source: Frerich RR, Aneshensel CS and Clark VA (1981) Prevalence of depression in Los Angeles County. *American Journal of Epidemiology* 113: 691-99.

Before you begin: Download from the course website: depress_small.Rdata

Background

The data for this illustration is a **subset of n=294** observations from the original study of 1000 adult residents of Los Angeles County. The purpose of the original study was to estimate the prevalence of depression and to identify the predictors of, and outcomes associated with, depression. The study design was a longitudinal one that included four interviews

In this illustration, only data from the first interview are used. Thus, this example is a cross-sectional analysis to identify the correlates of prevalent depression. Among these n=294, there are **50 events** of prevalent depression.

Codebook:

Variable	Label	Range/Codes
depressed	Case of depression	1=yes, 0 =no
age	Age, years	continuous
income	Income, thousands of dollars	continuous
female	Female gender	1=female, 0=male
unemployed	Unemployed	1=unemployed, 0=other
chronic	Chronic illness in past year	1=yes, 0=no
alcohol	Current alcohol use	1=yes, 0=no

Goal

Perform a multiple logistic regression analysis of these data to identify the correlates of prevalent depression.

R Illustration

```
# Load data. Keep only the variables of interest. Check.
load(file="depress_small.Rdata")
keepvars <- c("age", "alcohol", "chronic", "depressed", "female", "income", "unemployed")
temp <- depress_small[keepvars]
str(temp)
## 'data.frame': 294 obs. of 7 variables:
## $ age : num 68 58 45 50 33 24 58 22 47 30 ...
## $ alcohol : num 0 1 1 0 1 1 0 0 1 1 ...
## $ chronic : Factor w/ 2 levels "0. no", "1. yes": 2 2 1 2 1 2 2 1 2 1 ...
## $ depressed : num 0 0 0 0 0 0 0 0 1 0 ...
## $ female : num 1 0 1 1 1 0 1 0 1 0 ...
## $ income : num 4 15 28 9 35 11 11 9 23 35 ...
## $ unemployed: num 0 0 0 1 0 0 0 0 0 0 ...

# Descriptives: Continuous Variables - ALL
library(stargazer)
stargazer::stargazer(temp, type="text", median=TRUE, digits=4)
##
## =====
## Statistic N Mean St. Dev. Min Median Max
## -----
## age 294 44.4150 18.0854 18 42.5 89
## alcohol 294 0.7959 0.4037 0 1 1
## depressed 294 0.1701 0.3763 0 0 1
## female 294 0.6224 0.4856 0 1 1
## income 294 20.5748 15.2901 2 15 65
## unemployed 294 0.0476 0.2133 0 0 1
## -----
```

Because depressed is 0/1 mean=.17 says 17% of sample experienced event

```
# Descriptives: Continuous Variables - By Depression Status
depressed <- subset(temp, depressed==1)
paste("SUBSET: Depressed")
## [1] "SUBSET: Depressed"
stargazer::stargazer(depressed, type="text", median=TRUE, digits=4)
##
## =====
## Statistic N Mean St. Dev. Min Median Max
## -----
## age 50 40.3800 17.4003 18 34.5 79
## alcohol 50 0.8200 0.3881 0 1 1
## depressed 50 1.0000 0.0000 1 1 1
## female 50 0.8000 0.4041 0 1 1
## income 50 15.2000 9.8375 2 13 45
## unemployed 50 0.1200 0.3283 0 0 1
## -----

normal <- subset(temp, depressed==0)
paste("SUBSET: Normal")
## [1] "SUBSET: Normal"
stargazer::stargazer(normal, type="text", median=TRUE, digits=4)
##
## =====
## Statistic N Mean St. Dev. Min Median Max
## -----
## age 244 45.2418 18.1465 18 43.5 89
## alcohol 244 0.7910 0.4074 0 1 1
## depressed 244 0.0000 0.0000 0 0 0
## female 244 0.5861 0.4935 0 1 1
## income 244 21.6762 15.9755 2 17 65
## unemployed 244 0.0328 0.1784 0 0 1
## -----
```

```
# Descriptives: Discrete Variables - By Depression Status
library(summarytools)
```

```
# ALCOHOL
```

```
summarytools::cTable(temp$alcohol,temp$depressed,prop = 'r',totals=TRUE)
```

```
## Cross-Tabulation / Row Proportions
```

```
## Variables: alcohol * depressed
```

```
## Data Frame: temp
```

```
##
```

```
## -----
##      depressed      0      1      Total
##  alcohol
##      0           51 (85.00%)   9 (15.00%)   60 (100.00%)
##      1          193 (82.48%)  41 (17.52%)  234 (100.00%)
##      Total        244 (82.99%)   50 (17.01%)  294 (100.00%)
## -----
```

```
fisher.test(table(temp$alcohol,temp$depressed),conf.int=FALSE)
```

```
## Fisher's Exact Test for Count Data
```

```
##
```

```
## data: table(temp$alcohol, temp$depressed)
```

```
## p-value = 0.7049
```

The Null of "no association" is NOT rejected

```
## alternative hypothesis: true odds ratio is not equal to 1
```

```
## sample estimates:
```

```
## odds ratio
```

```
## 1.203052
```

Interpretation: In unadjusted analysis, depression is more prevalent among drinkers (Estimated OR = 1.2), but is not statistically significant (p-value = .70).

```
# CHRONIC
```

```
summarytools::cTable(temp$chronic,temp$depressed,prop = 'r',totals=TRUE)
```

```
## Cross-Tabulation / Row Proportions
```

```
## Variables: chronic * depressed
```

```
## Data Frame: temp
```

```
##
```

```
## -----
##      depressed      0      1      Total
##  chronic
##      0. no          126 (86.90%)   19 (13.10%)   145 (100.00%)
##      1. yes          118 (79.19%)   31 (20.81%)   149 (100.00%)
##      Total        244 (82.99%)   50 (17.01%)   294 (100.00%)
## -----
```

```
fisher.test(table(temp$chronic,temp$depressed),conf.int=FALSE)
```

```
## Fisher's Exact Test for Count Data
```

```
##
```

```
## data: table(temp$chronic, temp$depressed)
```

```
## p-value = 0.08876
```

The estimated association is at most marginally significant

```
## alternative hypothesis: true odds ratio is not equal to 1
```

```
## sample estimates:
```

```
## odds ratio
```

```
## 1.738905
```

Interpretation: In unadjusted analysis, depression is slightly more prevalent among the chronically ill (OR = 1.7), but only marginally statistically significant (p-value = .09).

```
# FEMALE at birth
summarytools::cTable(temp$female,temp$depressed,prop = 'r',totals=TRUE)
## Cross-Tabulation / Row Proportions
## Variables: female * depressed
## Data Frame: temp
##
## -----
##           depressed           0           1           Total
## female
##           0           101 (90.99%)   10 ( 9.01%)   111 (100.00%)
##           1           143 (78.14%)   40 (21.86%)   183 (100.00%)
##           Total           244 (82.99%)   50 (17.01%)   294 (100.00%)
## -----
```

```
fisher.test(table(temp$female,temp$depressed),conf.int=FALSE)
## Fisher's Exact Test for Count Data
##
## data: table(temp$female, temp$depressed)
## p-value = 0.004018
## alternative hypothesis: true odds ratio is not equal to 1
## sample estimates:
## odds ratio
## 2.815991
```

The null of "no association" is rejected

Interpretation: In unadjusted analysis, compared to males at birth, depression is more prevalent (Estimated OR = 2.8) among females at birth and is statistically significant (p-value = .004).

```
# UNEMPLOYED
summarytools::cTable(temp$unemployed,temp$depressed,prop = 'r',totals=TRUE)
## Cross-Tabulation / Row Proportions
## Variables: unemployed * depressed
## Data Frame: temp
##
## -----
##           depressed           0           1           Total
## unemployed
##           0           236 (84.29%)   44 (15.71%)   280 (100.00%)
##           1           8 (57.14%)    6 (42.86%)    14 (100.00%)
##           Total           244 (82.99%)   50 (17.01%)   294 (100.00%)
## -----
```

```
fisher.test(table(temp$unemployed,temp$depressed),conf.int=FALSE)
## Fisher's Exact Test for Count Data
##
## data: table(temp$unemployed, temp$depressed)
## p-value = 0.01825
## alternative hypothesis: true odds ratio is not equal to 1
## sample estimates:
## odds ratio
## 3.995901
```

The null of "no association" is rejected

Interpretation: Compared to employed persons, depression is more prevalent among the unemployed (OR = 4.0) and is statistically significant (p-value = .02).

Step 1. Fit single predictor models - Using Wald Z-score, retain predictors with significance levels $< .25$ or that are of a priori interest.

R

```
library(stargazer)
# Dear class: For now, I have commented out the detailed display of model summaries.

m_age <- glm(depressed ~ age, data=temp, family=binomial)
# summary(m_age)
m_alcohol <- glm(depressed ~ alcohol, data=temp, family=binomial)
# summary(m_alcohol)
m_chronic <- glm(depressed ~ chronic, data=temp, family=binomial)
# summary(m_chronic)
m_female <- glm(depressed ~ female, data=temp, family=binomial)
# summary(m_female)
m_income <- glm(depressed ~ income, data=temp, family=binomial)
# summary(m_income)
m_unemployed <- glm(depressed ~ unemployed, data=temp, family=binomial)
# summary(m_unemployed)

stargazer::stargazer(m_age, m_alcohol, m_chronic, m_female, m_income, m_unemployed, type="text", title="Single Predictor Models: beta and (SE[beta])", font.size="small", align=TRUE)
## Single Predictor Models: beta and (SE[beta])
## =====
##                               Dependent variable:
##                               -----
##                               depressed
##                               (1)      (2)      (3)      (4)      (5)      (6)
## -----
## age                        -0.016*
##                               (0.009)
##
## alcohol                     0.185
##                               (0.400)
##
## chronic1. yes               0.555*
##                               (0.318)
##
## female                     1.039***
##                               (0.377)
##
## income                     -0.036***
##                               (0.013)
##
## unemployed                 1.392**
##                               (0.564)
##
## Constant                   -0.917** -1.735*** -1.892*** -2.313*** -0.938*** -1.680***
##                               (0.404)  (0.362)  (0.246)  (0.332)  (0.266)  (0.164)
## -----
## Observations                294      294      294      294      294      294
## Log Likelihood             -132.514 -133.952 -132.504 -129.699 -129.701 -131.333
## Akaike Inf. Crit.         269.029  271.904  269.008  263.398  263.402  266.666
## =====
## Note:                       *p<0.1; **p<0.05; ***p<0.01
```

Interpretation: We will drop alcohol (p-value $> .10$) and consider further: age, chronic, female, income, unemployed

Step 2 – Assess candidate predictors for evidence of multicollinearity

R

```
library(car)
fit3 <- glm(depressed ~ age + chronic + female + income + unemployed, data=temp, family="binomial")
car::vif(fit3)
```

```
##      age      chronic      female      income unemployed
## 1.091355 1.114703 1.048511 1.055730 1.052486
```

Interpretation: Recall multicollinearity. Collinearity occurs when the predictors are themselves interrelated. If extreme this is a problem because a) the model is unstable; and b) the model is uninterpretable. Collinearity is suspected for VIF values > 10. Here things look okay so we will forge on!

Step 3. Fit multiple predictor model using step 1 predictors having crude significance < .25

R

```
fit4 <- glm(depressed ~ age + female + income + unemployed + chronic, data=temp, family=binomial)
summary(fit4)
```

```
## Call:
## glm(formula = depressed ~ age + female + income + unemployed +
##      chronic, family = binomial, data = temp)
##
```

```
## Deviance Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -1.2614  -0.6603  -0.4652  -0.3099   2.6026
```

```
## Coefficients:
```

```
##      Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -1.031844   0.612135 -1.686   0.0919 .
## age          -0.021938   0.009494 -2.311   0.0208 *
## female       0.812132   0.396880  2.046   0.0407 *
## income      -0.032067   0.014140 -2.268   0.0233 *
## unemployed   1.069739   0.598925  1.786   0.0741 .
## chronic1. yes 0.594859   0.350866  1.695   0.0900 .
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
## (Dispersion parameter for binomial family taken to be 1)
```

```
##      Null deviance: 268.12  on 293  degrees of freedom
## Residual deviance: 242.08  on 288  degrees of freedom
## AIC: 254.08
```

```
## Number of Fisher Scoring iterations: 5
```

Step 3 – Summary

Predictor	Adjusted Significance (Wald)	Remark
age	.021	Retain – pvalue is < .10
chronic	.090	For illustration purposes, let's consider dropping this variable, despite pvalue < .10 (it's close!)
female	.041	Retain – pvalue is < .10
income	.023	Retain – pvalue is < .10
unemployed	.074	Retain – pvalue is < .10.

Step 4. Fit the multivariable model containing predictors with adjusted significance levels < .10 from step 3. We will then compare the step 3 model with the step 4 model using a likelihood ratio test.

R

```
fit5 <- glm(depressed ~ age + female + income + unemployed,
            data=temp,
            family="binomial")

summary(fit5)
## Call:
## glm(formula = depressed ~ age + female + income + unemployed,
##      family = "binomial", data = temp)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.3604  -0.6459  -0.4943  -0.3173   2.5279
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.896828   0.597888  -1.500   0.1336
## age         -0.018802   0.009179  -2.048   0.0405 *
## female       0.938952   0.388746   2.415   0.0157 *
## income      -0.033431   0.014152  -2.362   0.0182 *
## unemployed   0.963457   0.592199   1.627   0.1038
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 268.12  on 293  degrees of freedom
## Residual deviance: 245.04  on 289  degrees of freedom
## AIC: 255.04
##
## Number of Fisher Scoring iterations: 5
```

Step 5. Investigate confounding.

Tentatively, a “good” final model is the four predictor model with predictors: **age**, **female**, **income**, and **unemployed**. Here, we explore possible confounding of the four predictor model by the omitted variable **chronic**. Specifically, we assess **chronic** as a potential confounder using 2 criteria:

- ___1. Likelihood Ratio test < .10 (or .05 or threshold of choice).
- ___2. Relative Change in estimated betas > 15% (or threshold of choice) using the following formula:

$$\Delta\hat{\beta} = \left(\frac{|\hat{\beta}_{\text{without confounder}} - \hat{\beta}_{\text{with confounder}}|}{\hat{\beta}_{\text{with confounder}}} \right) \times 100$$

By Hand: Likelihood ratio test comparing step 3 and step 4 models

$$\begin{aligned} \text{LR Test} &= [(-2) \ln (L)_{\text{REDUCED}}] - [(-2) \ln (L)_{\text{FULL}}] \\ &= [245.04] - [242.08] \\ &= 2.96 \end{aligned}$$

$$\text{LR Test df} = \Delta \text{Deviance df} = \Delta \# \text{ predictors in model} = 290 - 289 = 1$$

$$\text{p-value} = \text{Pr} \{ \text{Chi square with 1 degree of freedom} \geq 2.96 \} = .0853$$

This is not significant. We conclude that, in adjusted analysis, "chronic" is not statistically significant. Possibly, we can drop chronic

R

```
library(lmtest)
lmtest::lrtest(fit4, fit5)
## Likelihood ratio test
##
## Model 1: depressed ~ age + female + income + unemployed + chronic
## Model 2: depressed ~ age + female + income + unemployed
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    6 -121.04
## 2    5 -122.52 -1  2.9553    0.0856 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```


For Completeness Sake: Looking for $\geq 15\%$ Change in Betas for Predictors in Model

Potential confounding of **age, female, income, unemployed**

By: **chronic**

$\hat{\beta}_{\text{age}}(\text{w/o chronic}) = -.018802$; $\hat{\beta}_{\text{age}}(\text{w chronic}) = -.0219383$; Change = **14.30%**

$\hat{\beta}_{\text{female}}(\text{w/o chronic}) = .938952$; $\hat{\beta}_{\text{female}}(\text{w chronic}) = .8121316$; Change = **15.62%**

$\hat{\beta}_{\text{income}}(\text{w/o chronic}) = -.0334314$; $\hat{\beta}_{\text{income}}(\text{w chronic}) = -.0320672$; Change = **2.32%**

$\hat{\beta}_{\text{unemployed}}(\text{w/o chronic}) = .9634566$; $\hat{\beta}_{\text{unemployed}}(\text{w chronic}) = 1.069739$; Change = **9.94%**

The relative change in the beta for female is borderline at 15.6%. For parsimony, let's drop chronic.

Step 6. Investigate effect modification.

Are individuals who are both unemployed and with low income more likely to be depressed? For this illustration, we will create a new variable called **low** to capture individuals whose income is less than \$10,000. Then we will create an interaction of **low** and **unemployed**. **Tip** – When assessing interaction, it is necessary to include the main effects of both of the variables contributing to the interaction. Thus, this model includes the main effects **low** and **unemployed** in addition to the interaction **low_unemployed**.

R

```
library(lmtest)
fit6 <- glm(depressed ~ age + female + income + unemployed + low + low_unemployed,
            data=temp,
            family=binomial)

summary(fit6)
## Call:
## glm(formula = depressed ~ age + female + income + unemployed +
##     low + low_unemployed, family = binomial, data = temp)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.5506  -0.6523  -0.4798  -0.2754   2.6609
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.474687   0.683715  -0.694  0.48751
## age          -0.014759   0.009597  -1.538  0.12408
## female        1.036787   0.398427   2.602  0.00926 **
## income       -0.054349   0.020100  -2.704  0.00685 **
## unemployed    0.254521   0.875908   0.291  0.77137
## low          -0.945009   0.472267  -2.001  0.04539 *
## low_unemployed 1.544647   1.247602   1.238  0.21568
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 268.12  on 293  degrees of freedom
## Residual deviance: 240.38  on 287  degrees of freedom
## AIC: 254.38
##
## Number of Fisher Scoring iterations: 5
```

```
fit_temp <- glm(depressed ~ age + female + income + unemployed + low,
               data=temp,
               family=binomial)

lmtest::lrtest(fit6, fit_temp)
## Likelihood ratio test
##
## Model 1: depressed ~ age + female + income + unemployed + low + low_unemployed
## Model 2: depressed ~ age + female + income + unemployed + low
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    7 -120.19
## 2    6 -120.99 -1  1.6026    0.2055
```

Note - The lack of statistical significance is not surprising given the small number, 7, who are both UNEMPLOYED and with income < \$10,000 (R code and output for this is not shown - sorry). We will drop low and drop its interaction with unemployed.

Conclusion:

A reasonable multiple predictor model of depression in this sample contains the following predictors: age, female, income, and unemployed.

Examination of this model fit (see again the summary for fit5 on pages 33) suggests that, in adjusted analysis:

- (1) Older age is marginally associated with lower prevalence of depression.
Relative odds (OR) of depression associated with 1 year increase

$$= \exp(\beta_{AGE}) = \exp(-0.0188) = .98 \text{ (p=.04)}$$
- (2) Females, compared to males are more likely to be depressed.
Relative Odds (Odds ratio),

$$= \exp(\beta_{FEMALE}) = \exp(0.938952) = 2.6 \text{ (p=.016)}$$
- (3) Higher income is associated with lower prevalence of depression.
Relative odds (OR) of depression associated with \$1K increase

$$= \exp(\beta_{INCOME}) = \exp(-0.033431) = .97 \text{ (p=.018)}$$
- (4) Unemployed persons, are marginally significantly more likely to be depressed.
Relative Odds, OR

$$= \exp(\beta_{UNEMPLOYED}) = \exp(0.963457) = 2.6 \text{ (p=.010)}$$

8. Regression Diagnostics

With a fitted model come two applications, prediction and hypothesis tests.

- We've seen that a prediction is a guess of the expected outcome for a person with a particular profile of values of the explanatory variables (eg – value of vdtexpos) using the values of the estimated betas is obtained using the estimated betas:

$$\text{Predicted probability}_{\text{vdtexpos}} = \hat{\pi} = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1[\text{vdtexpos}])}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1[\text{vdtexpos}])}$$

- An example of an hypothesis test is the hypothesis test of the significance of VDTXPOS. The likelihood ratio test that the β for VDTXPOS is equal to zero compares

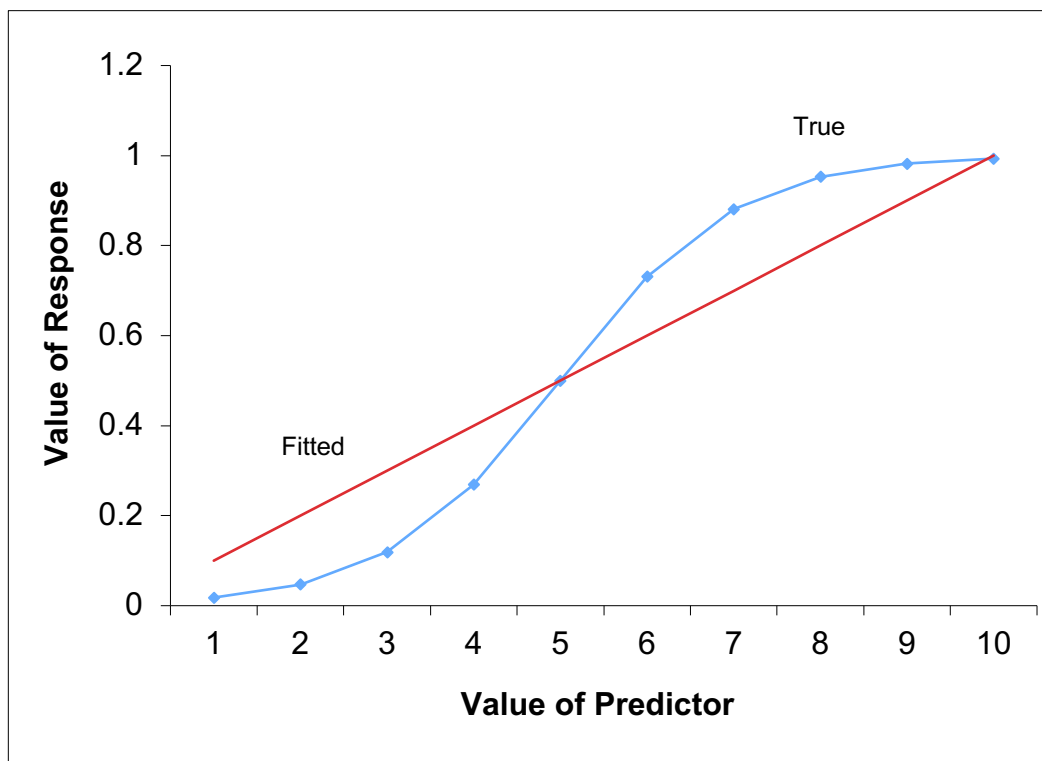
- the odds of SAB for exposed persons (“comparison”), versus
- the odds of SAB for Unexposed (“reference”) persons.

Neither prediction nor hypothesis tests have meaning when the model is a poor fit to the data.

Reasons for a poor fit include the following:

- The wrong relationship was fit.
- The data include extreme values which influence too greatly the fitted line.
- Important explanatory variables have not been included.

We need **regression diagnostics** for the detection of a poor fit:



Example - The fit is poor here because the true relationship is quadratic, not linear.

We notice that the discrepancies between the observed and the fitted values are not of consistent size.

Some are large and some are small.

Goodness-of-fit assessments are formal techniques for identifying such inconsistencies.

These techniques become especially important when a picture is not possible, as when the number of predictors is greater than one.

Assessing regression model adequacy was introduced previously (Unit 5, Regression and Correlation). Regression diagnostics are of two types:

- Systematic component
 - Is the assumption of linearity on the $\ln(\text{odds})$ scale correct?
 - Is the logistic model formulation a reasonably good fit?
 - Should we have fit a different model?
 - Does the fitted model predict well?
- Case analysis
 - Is the fitted model excessively influenced by one or a small number of individuals?

There exist methods to address each of these regression diagnostic questions.

Question	Method of Assessment
Is the assumption of linearity on the $\ln(\text{odds})$ scale correct?	a. Assessment of linearity
Is the logistic model formulation a reasonably good fit?	b. Hosmer-Lemeshow test for overall goodness of fit.
Should we have fit a different model?	c. Linktest
Does the fitted model predict well?	d. Classification table e. The ROC Curve
Is the fitted model excessively influenced by one or a small number of individuals or <u>covariate patterns</u> ? <i>Note – Here we might look at covariate patterns instead of individuals.</i>	f. Pregibon Delta beta statistic

a. Assessment of Linearity

A logistic regression model assumes that the **logit of the probability (π) of event occurrence** (eg – spontaneous abortion) is **linear** in the predictors X_1, X_2, \dots etc.

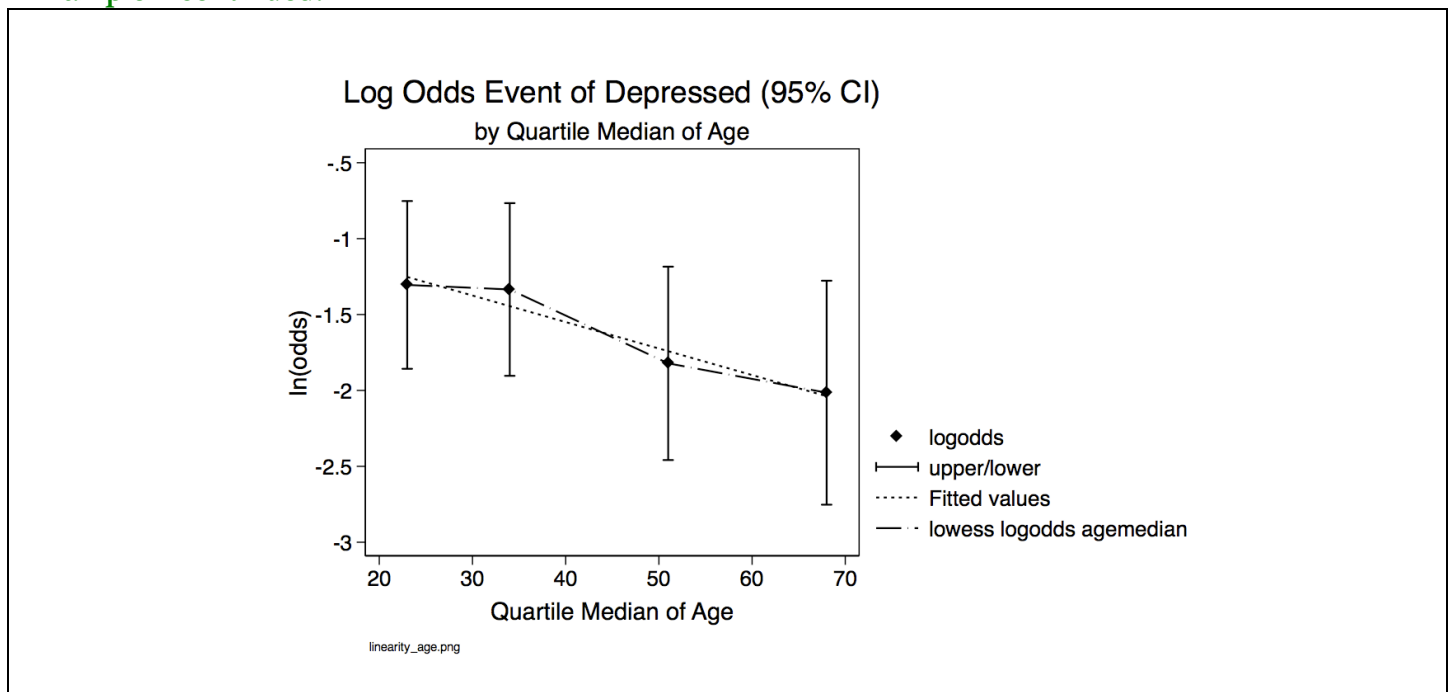
$$\text{logit}[\pi_x] = \text{logit}[E(Y)] = \ln\left[\frac{\pi_x}{1-\pi_x}\right] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_5 X_5 + \beta_6 X_6$$

Violation of the assumption of linearity of the logit in a continuous predictor can lead to incorrect estimates and incorrect conclusions. A variety of approaches are available for assessing the assumption of linearity in logistic regression but are *beyond the scope of these notes*.

A **graphical** assessment of linearity of $Y = \text{logit}$ with changes in $X = \text{predictor}$ involves five steps

1. Collapse the predictor values of X into groups (eg; quartiles)
2. In each group, obtain the median value of the predictor variable X .
3. In each group, obtain the observed proportion experiencing the event Y .
4. In each group, obtain the observed logit [proportion experiencing event]
Tip – Obtain 95% CI limits as well.
5. Produce a two-way plot of $X = \text{midpoint}$ versus $Y = \text{logit}$, perhaps with some overlays.

Example – continued.



Not bad! The plot looks reasonable enough that it is okay to model the logit linearly in age.

b. The Hosmer-Lemeshow Test of Goodness-of-Fit

The **Hosmer-Lemeshow Goodness of Fit Test** compares observed versus predicted counts of outcome events in each of several “meaningful” subgroups of the data, in a manner similar to the Chi Square Goodness of Fit Test introduced in Unit 4, Categorical Data. If the fit is good (null hypothesis is true), the observed and (model based) expected counts will be close and their differences will be small. The actual test statistic is a sum of $(\text{observed} - \text{expected})/\text{expected}^2$ and is distributed chi square under the null hypothesis.

Null Hypothesis: “Good fit” is indicated by similar counts of observed and predicted counts in all the subgroups.

The difference between the two counts is then close to zero.

The sum, taken over the subgroups, is also small.

The Groups Used in a Hosmer-Lemeshow Test are defined by the predicted probabilities

Within each group, members have similar predicted probabilities of outcome event.

The most commonly used groups are 10 subgroups defined by deciles of predicted.

1st subgroup: This is the 1/10th of sample of persons who have the **lowest predicted probabilities** of outcome event.

2nd subgroup: This is the next 1/10 of sample of persons. These persons have the **next lowest predicted probabilities** of outcome event.

And so on

10th subgroup: This is the last 1/10 of sample of persons. These persons have the **highest predicted probabilities** of outcome event.

Hosmer-Lemeshow Goodness of Fit Test

Yet another chi square test

H_O: The current model is a “good” fit to the data.

H_A: not.

$$\chi^2_{\text{Hosmer-Lemeshow; DF}=\# \text{ groups}-2} = \sum_{\text{decile of risk}} \left\{ \frac{[\text{Observed count} - \text{Predicted count}]^2}{\text{Predicted count}} \right\}$$

Rejection occurs for large values of the chi square statistic with associated small p-values

Calculation of **observed** and (model fit) **predicted** counts:

Observed count = Actual number of events in decile

Predicted count = (# in group) (Average predicted probability)

When the null hypothesis of a “good” fit is true,

$\chi^2_{\text{Hosmer-Lemeshow}}$ is distributed Chi Square, approximately. With df= (# groups) – (2)

For example, with 8 groups, the degrees of freedom = 6

Large values of this statistic suggest a poor fit.

Statistically significant values of the Hosmer-Lemeshow statistic evidence **ONLY** that the fit is poor. We do not learn why. Further assessments are necessary to understand their nature.

R

```
library(ResourceSelection)
hl <- ResourceSelection::hoslem.test(fit5$y, fitted(fit5), g=8)
hl
##
## Hosmer and Lemeshow goodness of fit (GOF) test
##
## data: fit5$y, fitted(fit5)
## X-squared = 0.96101, df = 6, p-value = 0.987

hl.test <- cbind(round(hl$observed,digits = 0),round(hl$expected,digits = 1))
colnames(hl.test) <- c("Obs_0", "Obs_1", "Exp_0", "Exp_1")
paste("Hosmer-Lemeshow Goodness of Fit Test of 'final' Model")
## [1] "Hosmer-Lemeshow Goodness of Fit Test of 'final' Model"
hl.test

##              Obs_0 Obs_1 Exp_0 Exp_1              Obs_0 and Exp_0 counts look similar (nice)
## [0.0143,0.0624]    35     2  35.5   1.5              Obs_1 and Exp_1 counts look similar as well (it follows)
## [0.0624,0.0807]    35     2  34.4   2.6
## [0.0807,0.118]     32     4  32.4   3.6
## [0.118,0.158]      32     5  31.8   5.2
## [0.158,0.18]       32     5  30.7   6.3
## [0.18,0.223]       29     7  28.7   7.3
## [0.223,0.302]      26    11  27.4   9.6
## [0.302,0.646]      23    14  23.1  13.9
```

The null hypothesis is NOT rejected. Good news; we don't want to have to worry about a bad fit. The Hosmer_Lemeshow test ($p=0.987$) suggests no statistically significant departure from a good fit.

c. The Linktest

The **Link Test** is an example of a **specification test**.

Like the Hosmer-Lemeshow statistic, the **Link Test** is a simple check of the fitted model. It assesses whether or not the fitted model is adequate fit (null hypothesis) to the data or, if not, if there is still some additional modeling that needs to be done (alternative hypothesis). The crudeness of the Link Test is that what we learn is limited. If the null hypothesis is rejected, we know only that some alternative modeling is needed, but we don't know what alternative modeling is needed.

Link Test

H_0 : The current model is an adequate fit to the data.

H_A : Alternative modeling is needed.

A Likelihood Ratio (LR) Test is performed and compares a “null hypothesis” adequate model (reduced) with an “alternative hypothesis enhanced (full) model:

Reduced: $\text{logit}[\pi] = \beta_0 + \beta_1[\hat{\pi}_{\text{model}}]$

Full: $\text{logit}[\pi] = \beta_0 + \beta_1[\hat{\pi}_{\text{model}}] + \beta_2[\hat{\pi}_{\text{model}}^2]$

Thus,

H_0 : $\beta_2 = 0$

H_A : not

Key -

$\hat{\pi}_{\text{model}}$: This is the predicted probability from our model; we hope this is significant.

$\hat{\pi}_{\text{model}}^2$: If the null is true (the model is adequate), this should be **non-significant**.

Rejection of the null occurs for large values of the LR Test and associated small p-values.

R

```

hat <- predict(fit5)
hatsq <- hat^2
linktest <- summary(glm(depressed ~ hat + hatsq,data=temp, family=binomial))
paste("Link Test of Final Model")
## [1] "Link Test of Final Model"

paste("Null: Look for hatsq NOT significant (suggests all is well)")
## [1] "Null: Look for hatsq NOT significant (suggests all is well)"

linktest
##
## Call:
## glm(formula = depressed ~ hat + hatsq, family = binomial, data = temp)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.3961  -0.6421  -0.4929  -0.3236   2.5069
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.04389    0.50704   0.087   0.931
## hat         1.07581    0.65695   1.638   0.102
## hatsq       0.02519    0.20412   0.123   0.902
##                                     The null of "extra modeling needed" is NOT rejected
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 268.12  on 293  degrees of freedom
## Residual deviance: 245.02  on 291  degrees of freedom
## AIC: 251.02
##
## Number of Fisher Scoring iterations: 5

```

$\hat{\pi}_{\text{model}}$: This is marginally significant ($p=.10$); perhaps we'd hoped for better. But okay.

$\hat{\pi}_{\text{model}}^2$: This is **non-significant** ($p=.90$). Good news.

The Link Test ($p=.902$) suggests no statistically significant departure from model adequacy.
The null hypothesis of "model adequacy" is NOT rejected. **Good news!**

d. The Classification Table

Rationale

- Just because the fitted model is a good fit overall doesn't mean that individual predictions are correct most of the time.
- The classification table, and associated plots, are useful in a selected analysis setting:

We might want to use the fitted equation to make predictions about event occurrence, depending on an individual's profile of values on the variables in the model.

Method

- For each individual, there are two quantities to work with
 - Actual outcome: Yes/No indicator of event occurrence
 - Estimated probability of event: Between 0 and 1
- Choose a threshold probability for event declaration by model.
 - Default is usually 0.5
 - This can be reset.
 - Consideration of several permits construction of ROC curve.

A separate classification table is produced for each cut-off you select

		Observed (True)		
		Event	Non-Event	
<u>Predicted</u>	Event			
	Non-Event			

Example:

Suppose that for subject id=103 observed event = YES predicted probability = .68

When cut-off=.60 observed event is still = YES Now, predicted event = YES Because .68 > .60
 When cut-off=.70 observed event is still = YES But, predicted event = NO Because .68 < .70

R

```
mylogit <- fit5
# The function step( ) provides a tabulation of AIC values for various models given a set of predictors
mysteps <- step(mylogit, depressed ~ age + female + income + unemployed, data=temp, family="binomial")
## Start: AIC=255.04
## depressed ~ age + female + income + unemployed
##
##           Df Deviance   AIC
## <none>          245.04 255.04      This model contains all the predictors. AIC = 255.04
## - unemployed  1   247.54 255.54      This model omits the predictor "unemployed"
## - age         1   249.40 257.39      and so on....
## - female      1   251.56 259.56
## - income      1   251.73 259.73

classDF <- data.frame(response = temp$depressed, predicted = ifelse(fitted(mysteps)>=0.5,1,0))
paste("Classification Table: Cut-off=0.5")
## [1] "Classification Table: Cut-off=0.5"
xtabs(~ predicted + response, data = classDF)
##           response
## predicted  0    1
##           0 243  48
##           1   1   2
```

Using cutoff = .50, the # observations with perfect concordance = 243 + 2

```
classDF2 <- data.frame(response = temp$depressed, predicted = ifelse(fitted(mysteps)>=0.6,1,0))
classDF3 <- data.frame(response = temp$depressed, predicted = ifelse(fitted(mysteps)>=0.1,1,0))

paste("Classification Table: Cut-off=0.6")
## [1] "Classification Table: Cut-off=0.6"
xtabs(~ predicted + response, data = classDF2)
##           response
## predicted  0    1
##           0 243  49
##           1   1   1

paste("Classification Table: Cut-off=0.1")
## [1] "Classification Table: Cut-off=0.1"
xtabs(~ predicted + response, data = classDF3)
##           response
## predicted  0    1
##           0  84   7
##           1 160  43
```

Not surprising. Using cutoff = .50, the concordance is reduced = 84 + 43

e. The ROC Curve

One of the uses of a fitted logistic model is to make predictions for new individuals; eg – **is this new person predicted to experience the event or not?**

An ROC curve (“Receiver-Operating Characteristic”) is a visual display of the overall performance of a fitted logistic model and its associated equation for predicted probabilities. It takes into consideration that there are **two kinds of errors of prediction**: (1) a true event is predicted to be a non-event (false negative) and (2) a true non-event is predicted to be an event (false positive, which is the same as $1 - \text{specificity}$).

For various choices of “cut-off” (**recall - this is the value above which a predicted probability is classified as a predicted event**) an ROC curve is plot of $X = \text{false positive}$ against $Y = \text{true positive}$ values for various choices of “cut-off”:

“Cutoff”	.10	.20	etc	.80	.90
$X = \text{false positive} = 1 - \text{specificity}$					
$Y = \text{correct positive} = \text{sensitivity}$					

Key

- In a real world application, **the choice of “cutoff” has real world implications** as when a predicted event=yes prompts the initiation of treatment.
- A diagonal line with slope=1 is a reference line.** It represents the ROC curve for test that performs no better than the **flip of a coin**.
- The area under the ROC curve is often denoted c-statistic. It has a defined meaning:

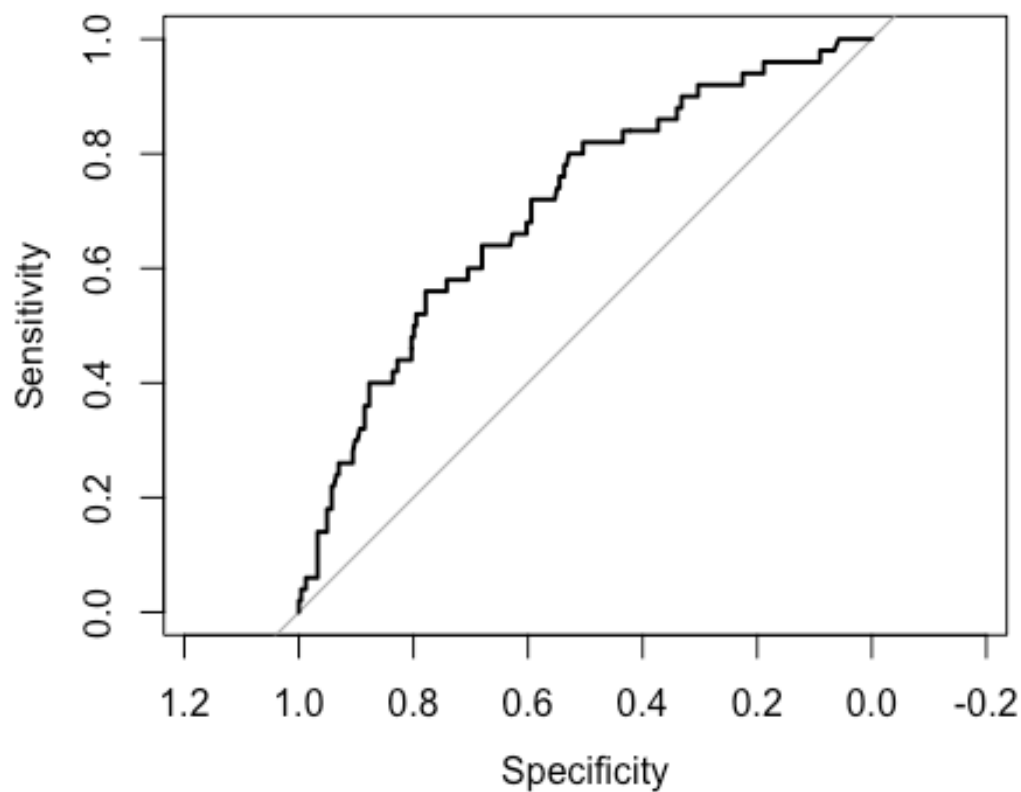
ROC Curve

c-statistic = Overall % correctly classified

= area under the curve

R

```
library(pROC)
pROC::roc(temp$depressed, fit5$fitted.values, plot=TRUE)
```



```
## Call:
## roc.default(response = temp$depressed, predictor = fit5$fitted.values, plot = TRUE)
##
## Data: fit5$fitted.values in 244 controls (temp$depressed 0) < 50 cases (temp$depressed 1).
## Area under the curve: 0.708
```

Key -

- Recall - The straight line with slope = 1 is a reference line; it corresponds to the ROC curve where chance alone is operating (coin toss with probability heads = .50)
- ROC c-statistic = .7080** says that the overall % who are correctly classified is 70.8%. This is not very impressive, actually. We typically hope to do better.

f. The Pregibon Delta Beta Statistic

Recall the **Cook's Distance Statistic** introduced in unit 5, **Regression and Correlation**. This statistic provides a measure of the extent to which inclusion or non-inclusion of an individual changes the estimated betas.

The plot is of $X = \text{Subject ID}$ versus $Y = \text{Cook's Distance}$
Spikes in the plot identify individuals whose inclusion are influential on the fit.

The analogue in logistic regression is the **Pregibon Delta Beta Statistic, dbeta**. The formula is beyond the scope of this course. However, a feel for it is the following:

dbeta = function of { standardized difference in betas w deletion of individual
or deletion of covariate pattern }

The Pregibon Delta Beta Statistic can be computed for **study individuals** or for **covariate patterns** instead of study id.

- A **covariate pattern** is a unique profile (or combination) of values on the variables.
- The **maximum number** of covariate patterns in a data set occurs when every individual is unique in his/her pattern of values of the predictors. In this extreme case, the **number of covariate patterns = sample size = n**.
- Often, however, the same covariate pattern is shared by more than one individual (eg – 4 subjects have age=50, sex=male, exposure=yes). Thus, often, the **number of covariate patterns < n**.

The plot is of $X = \text{predicted probability}$ versus $Y = \text{dbeta}$

- Small values of dbeta: individual or covariate pattern is not influential
Small: dbeta values less than 1 or so, approx
- Large values of dbeta: individual or covariate pattern is influential
Large: dbeta values > 1

Tip – Regardless of the magnitudes of the dbeta, be on the look out for spikes
Spikes are suggestive of comparative influence

9. Example - Disabling Knee Injuries in the US Army

Source: Sulsky SI, et al . Risk Factors for Disability Discharge from the US Army Related to Occupational Knee Injury (2000).

Background:

The strongest correlate of lost time from work, lost productivity, and lost working years of life is occupational injuries.

Occupational activities have been found to be associated with knee disorders.

Poorly understood, however, are the differences in risk of knee disorders associated with socio-demographic versus occupational task characteristics.

Better understanding of the socio-demographic variations in risk of occupational knee injury is important to future studies of occupational risks.

Therefore, Sulsky et al conducted a case-control study to investigate selected socio-demographic risk factors for occupational knee injury in the US Army.

Research Question:

What are the separate and joint effects of gender, age, and race/ethnicity in the odds of disabling knee injury among enlisted Army personnel on active duty between 1980 and 1994?

Design: Nested case-control investigation of knee related disability within the occupational cohort of enlisted US Army personnel on active duty between 1980 and 1994.

Total Army Injury and Health Outcomes Data Base (TAIHOD)

2.1 million males at birth
283,000 females at birth
 ≈ 2.4 million



Data Library	
Cases	Controls
First record of any of 11 eligible codes 7868 males at birth 860 females at birth 8728 total	Density sampling* of TAIHOD by year, separately for each gender 11,758 males at birth (control:case = 1.5:1) <u>5,109 females at birth (control:case = 6:1)</u> 16,867 Total (control:case = 2:1)



Analysis Sample			
	Cases	Controls	Control:Case
Females at birth	860: all cases	2580: density sampling by year*	3:1
Males at birth	1005: equal random sampling by year over 15 years (67/year)	3009: equal random sampling by year over 15 years (201/year)	3:1
Total	1865	5589	7454

* For the unfamiliar - Density Sampling by Year: For each year, controls were drawn in proportion to the number of cases for that year. (E.g. – A year with 2 cases and 3:1 sampling of controls yields 6 controls for that year.)

Estimated Distribution of Risk Factors: Age and Race/Ethnicity, by Sex at Birth

Our estimates will have to take into account the method of sampling employed. How does this work?

Let's look at a simple illustration. Suppose

Males at birth	Females at birth
Source Population, N=2000 Size of random sample, n=100 Probability[inclusion] = $100/2000 = .05$ Weight per person included = $1/.05 = 20$ Each male at birth in the sample represents 20 males at birth in the source population.	Source Population, N=1000 Size of random sample, n=100 Probability[inclusion] = $100/1000 = .10$ Weight per person included = $1/.10 = 10$ Each female at birth in the sample represents 10 females at birth in the source population.
The number of males at birth <21 years of age in the <u>sample</u> is # = 50. Therefore, <u>estimated</u> number of males at birth <21 years of age in the source <u>population</u> is $50 \times (\text{weight}=20) = 1000$	The number of females at birth <21 years of age in the <u>sample</u> is # = 25 Therefore, <u>estimated</u> number of females at birth <21 years of age in the source <u>population</u> is $25 \times (\text{weight}=10) = 250$

What is the overall relative frequency of age < 21 years?

Unweighted estimate describes the sample: $(50+25)/200 = 37.5\%$.
Weighted estimate describes the population: $= (1000+250)/3000 = 41.7\%$

REMINDER

When a study calls for stratified sampling with disproportionate sampling of selected groups, estimates of population characteristics must take sample weights and stratified sampling into account.

**Estimated Distribution of Risk Factors:
Age and Race/Ethnicity, by Sex at birth**

			Relative Frequency* Among	
			<u>Cases</u>	<u>Controls</u>
Males at birth	Age	<21	15	20
		21-23	19	19
		23-26	26	20
		26-30.36	20	18
		30.36-54	19	23
	Race/Ethnicity	Unknown	0	0
		White	71	62
		Black	22	29
		Other	7	9
Females at birth	Age	<21	19	19
		21-23	18	20
		23-26	19	22
		26-30.36	24	23
		30.36-54	20	16
	Race/Ethnicity	Unknown	0.2	0
		White	68	47
		Black	26	45
		Other	6	8

- Estimated relative frequencies take sample weights and stratified sampling into account.

We'll use quintiles of age.

Race/Ethnicity will be categorized as White/Non-White.

A multivariable logistic regression model analysis will explore the separate and joint associations with disabling knee injury of age, sex at birth, and race/ethnicity.

Recall the Research Question:

What are the separate and joint effects of sex at birth, age, and race/ethnicity in the odds of disabling knee injury among enlisted Army personnel on active duty between 1980 and 1994?

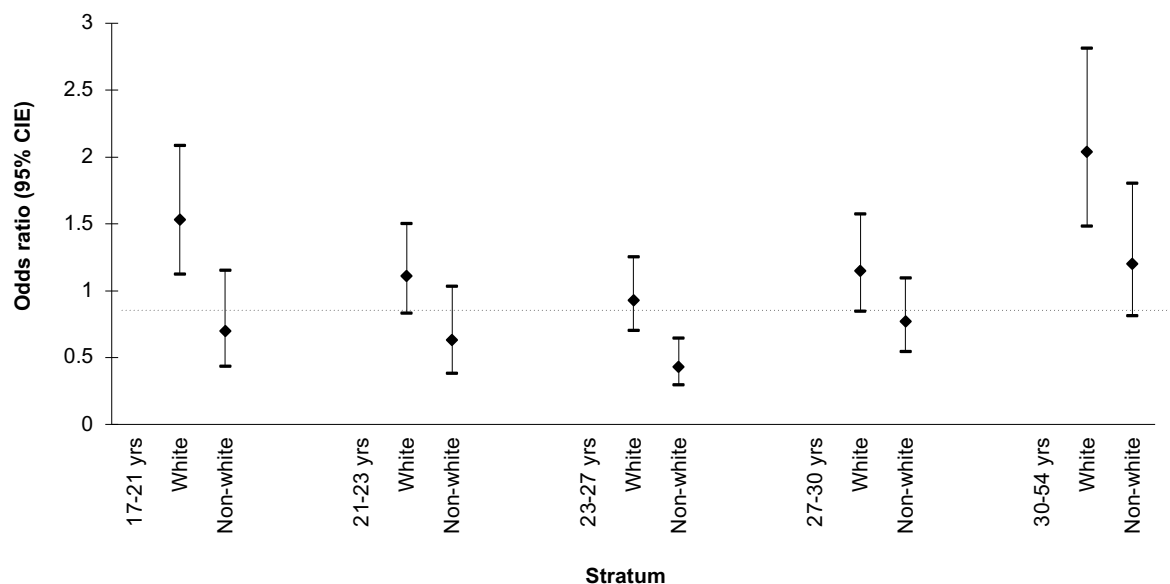
We are especially interested in identifying possible interactions.

- This analysis is to guide future analyses of occupational risk factors.
- A “traditional” analysis of occupational risk factors might simply control for age, sex at birth, and race/ethnicity.
- If interactions exist among age, sex at birth, and race/ethnicity, inclusion of only main effects might lead to incorrect inferences.

Therefore, the analysis plan seeks to estimate

- The separate effects of sex at birth on risk of disabling knee injury among groups defined by age | and race/ethnicity.
e.g. – Is the effect of sex at birth different among young workers compared to the effect of sex at birth among older workers?
- The separate effects of increasing age on risk of disabling knee injury among groups defined by sex at birth and race/ethnicity.
e.g. – Is the effect of increasing age different among males at birth and females at birth?

Figure 1: Relative odds of discharge for disabling knee injury among enlisted women compared to men, stratified by age (quintiles) and race.



- **Among Whites:**

Females at birth are at higher risk of disabling knee injury than males at birth at all ages except among persons aged 23-27.

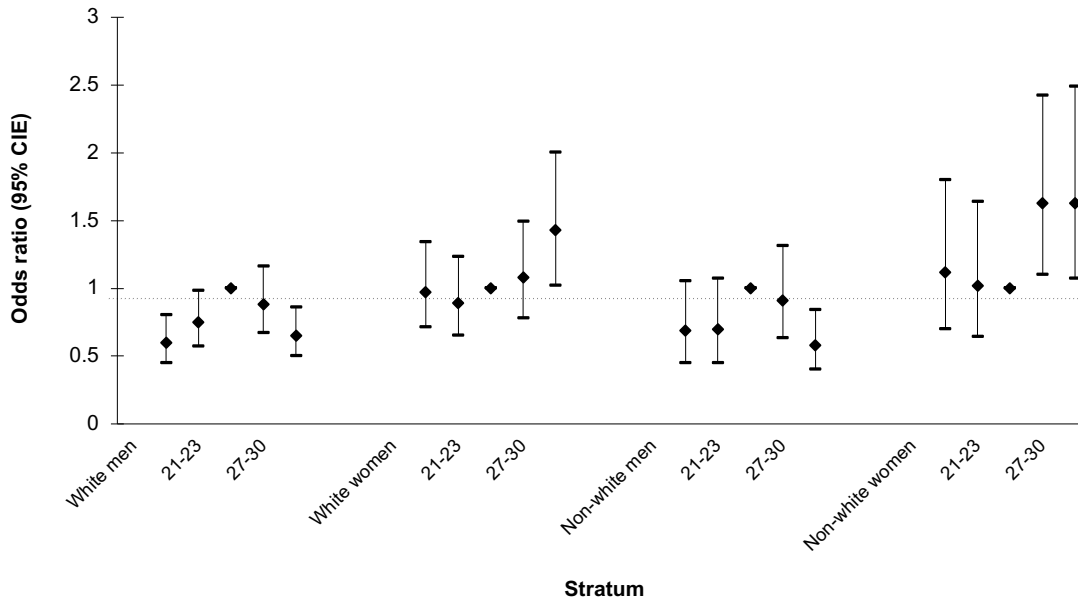
The sex at birth effect is greatest among the youngest (17-21 years) and oldest (30-54) persons. (“U” shape)

- **Among non-Whites:**

Females at birth are at lower risk of disabling knee injury than males at birth at all ages except among persons aged 30-54.

The sex at birth effect is greatest among persons in the middle age group (23-27 years). (“U” shape)

Figure 2: Relative odds of discharge for disabling knee injury with increasing age, stratified by sex and race.



note: The reference age group is age 23-27 years.

- **Among Males at birth:**

With increasing age, the change in risk of disabling knee injury exhibits a “ \cap ” pattern.

The “ \cap ” pattern among Whites is stronger than the “ \cap ” pattern among non-Whites.

- **Among Females at birth:**

With increasing age, the change in risk of disabling knee injury exhibits a “ \cup ” pattern.

The “ \cup ” pattern among Whites is more precise than the “ \cup ” pattern among non-Whites.

This example is a nice illustration of the distinction between confounding and effect modification

CAUTION!!

Confounding and effect modification are not simply about sampling and variations in nature. Their identification in statistical analysis is also a function of the choice of scale of measurement.

In the analysis of the relative odds of disabling knee injury, we are actually speaking of

Odds ratio confounding
Odds ratio modification

A (odds ratio) relationship between “E” and “D” that is confounded by X means:

- 1) X is related to both “E” and “D”
- 2) The unadjusted association between “E” and “D” is spuriously large or small because of the confounding effects of X
- 3) However, at each level of X, the association between “E” and “D” is the same.
- 4) A logistic regression analysis of the “E”-“D” relationship should include the predictor variable X.

A (odds ratio) relationship between “E” and “D” that is modified by X means:

- 1) X is related to both “E” and “D”
- 2) With changes in the level of X, the association between “E” and “D” changes also.
- 3) A logistic regression analysis of the “E”-“D” relationship should reveal these changes with X through the inclusion of “E”-“X” interactions.

Appendix Overview of Maximum Likelihood Estimation

The method of maximum likelihood estimation is used to obtain “good” guesses of the values of the regression coefficients, $\beta_0 \dots \beta_6$.

What do we mean by “good”?

1) Recall that, in linear model regression, “good” was conceptualized as obtaining guesses of $\beta_0 \dots \beta_6$ that make as small as possible the total of the vertical distances between the observed data Y and the fitted values \hat{Y} . We use the method of least squares and choose guesses, represented as $\hat{\beta}_0 \dots \hat{\beta}_6$, which minimize the residual sum of squares:

$$\text{Residual sum of squares} = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^N \left(Y_i - [\hat{\beta}_0 + \dots + \hat{\beta}_6 x_i] \right)^2$$

When the distribution of the errors is normal, we have a very nice result:

Method of least squares = Method of maximum likelihood; where

“maximum likelihood estimation” is described below.

2) In logistic model regression, “good” is obtaining guesses of $\beta_0 \dots \beta_6$ which make as large as possible the likelihood of obtaining the observed data. This is the method of maximum likelihood.

A Feel for Maximum Likelihood Estimation

A box contains two coins, A and B. One is selected.

“A” is fair and lands “heads” with probability $\pi = .50$.

“B” is not fair. It lands “heads” with probability $\pi = .67$.

Game: Toss the coin $n=20$ times. Note how many times the coin lands “heads”. Call this X .
Suppose $X=15$.

Question: Which choice of π , .50 or .67, maximizes the chances that the coin lands “heads” 15 times?

$\binom{20}{15} \pi^{15} (1-\pi)^{20-15}$	$\pi = .50$	$\pi = .67$
Likelihood, L L = Prob [X=15]	= .10	= .45

Review: The expression $\binom{20}{15}$ is a binomial coefficient and represents the number of ways to choose 15 items from 20. It is equal to $20!/[15! 5!]$.

There is a 10% chance of 15 “heads” when $\pi = .50$. There is a 45% chance of 15 “heads” when $\pi = .67$.

Even though scenarios of low probability do occur, the maximum likelihood estimate of the unknown probability of heads is chosen to be the one that makes as large as possible, the likelihood of the actual data.

\Rightarrow The maximum likelihood guess of $\pi = .67$.

Overview of Maximum Likelihood Estimation in Logistic Regression

Preliminaries

- (1) It is assumed that the n outcomes Y_1, \dots, Y_n are independent
- (2) It is also assumed that each Y_i is the outcome of a Bernoulli (π_i) trial
- (3) We'll use the notation L_i to represent each individual "likelihood", also called the probability density:

$$\begin{aligned} L_i &= \text{Probability}[Y_i=y_i] \\ &= \pi_i^{y_i} (1-\pi_i)^{1-y_i} \\ &= \left[\frac{\pi_i}{1-\pi_i} \right]^{y_i} (1-\pi_i)^1 \end{aligned}$$

- (4) We'll use the notation L to represent the likelihood of all n observations in the data. This is also called the "probability density of the data"

L = likelihood of the data

$$\begin{aligned} L &= \text{Probability}[Y_1=y_1, Y_2=y_2, \dots, Y_p=y_p] \\ &= \text{Probability}[Y_1=y_1] \text{Probability}[Y_2=y_2] \dots \text{Probability}[Y_p=y_p] \text{ by independence} \\ &= \prod_{i=1}^n \text{Probability}[Y_i=y_i] \\ &= \prod_{i=1}^n L_i \end{aligned}$$

- (4) The logistic model with predictors $\beta_0, \beta_1, \dots, \beta_p$ is defined

$$\ln \left[\frac{\pi_i}{1-\pi_i} \right] = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi}$$



x_{1i} = value of the variable X_1 for the "ith" person, etc.

(5) The logistic model with predictors $\beta_0, \beta_1, \dots, \beta_p$ also means that

$$\begin{aligned}\ln(1-\pi_x) &= \ln \left[\frac{1}{1+\exp(\beta_0+\beta_1x_1+\dots+\beta_px_p)} \right] \\ &= \ln[1] - \ln[1+\exp(\beta_0+\beta_1x_1+\dots+\beta_px_p)] \text{ because } \ln(a/b) = \ln(a) - \ln(b) \\ &= 0 - \ln[1+\exp(\beta_0+\beta_1x_1+\dots+\beta_px_p)] \text{ because } \ln[1]=0 \\ &= -\ln[1+\exp(\beta_0+\beta_1x_1+\dots+\beta_px_p)]\end{aligned}$$

Overview

- Maximum likelihood estimation of $\beta_0, \beta_1, \dots, \beta_p$ is accomplished by maximizing the natural logarithm of the likelihood L of the data.
- We'll let $L(\beta) = \ln \{ L \}$ represent the natural logarithm of the data under the assumption of the logistic regression model.

Solution for $L(\beta)$.

This is the function of the data that we seek to maximize with respect to $\beta_0, \beta_1, \dots, \beta_p$

$$L(\beta) = \ln \{ L \}$$

$$\begin{aligned} &= \ln \left[\prod_{i=1}^n L_i \right] \\ &= \sum_{i=1}^n \{ \ln[L_i] \} \quad \text{because } \ln[(a)(b)] = \ln(a) + \ln(b) \\ &= \sum_{i=1}^n \ln \left\{ \left[\frac{\pi_i}{1-\pi_i} \right]^{y_i} (1-\pi_i)^l \right\} \quad \text{by preliminary \#3} \\ &= \sum_{i=1}^n \left\{ \ln \left[\frac{\pi_i}{1-\pi_i} \right]^{y_i} + \ln(1-\pi_i)^l \right\} \quad \text{again because } \ln[(a)(b)] = \ln(a) + \ln(b) \\ &= \sum_{i=1}^n \left\{ y_i \ln \left[\frac{\pi_i}{1-\pi_i} \right] + \ln(1-\pi_i) \right\} \quad \text{because } \ln(a^b) = (b) \ln[a] \\ &= \sum_{i=1}^n \left\{ y_i \ln \left[\frac{\pi_i}{1-\pi_i} \right] \right\} + \sum_{i=1}^n \{ \ln(1-\pi_i) \} \\ &= \sum_{i=1}^n \left\{ y_i \left[\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi} \right] \right\} + \sum_{i=1}^n \{ \ln(1-\pi_i) \} \quad \text{by preliminary \#4} \\ &= \sum_{i=1}^n \left\{ y_i \left[\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi} \right] \right\} - \sum_{i=1}^n \left\{ \ln \left(1 + \exp \left[\beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi} \right] \right) \right\} \quad \text{by preliminary \#5} \end{aligned}$$

Maximization of the Log-Likelihood $L(\beta) = \ln \{ L \}$

Maximizing $L(\beta) = \ln \{ L \}$ with respect to each of $\beta_0, \beta_1, \dots, \beta_p$ is not the straightforward solution that was seen for estimating β_0 and β_1 in simple linear regression. It is beyond the scope of this course to develop the solution required here.

In brief, the solution for the maximum likelihood estimates is obtained by a method called [Newton Raphson iteration](#). In brief, this iterative procedure for maximizing $L(\beta) = \ln \{ L \}$ works with a linear approximation of the derivative of $L(\beta) = \ln \{ L \}$ with respect to $\beta_0, \beta_1, \dots, \beta_p$ and an initial estimate of $\beta_0, \beta_1, \dots, \beta_p$. From there an updated estimate of $\beta_0, \beta_1, \dots, \beta_p$ is obtained. Iteration continues until a convergence criterion is reached.