

# **Binary Response: Logistic Regression**

STAT 526  
Professor Olga Vitek

March 29, 2011

# **Model Specification and Interpretation**

# Probability Distribution of a Binary Outcome $Y$

- In many situations, the response variable has only two possible outcomes
  - Disease ( $Y = 1$ ) vs Not diseased ( $Y = 0$ )
  - Employed ( $Y = 1$ ) vs Unemployed ( $Y = 0$ )
- Response is *binary or dichotomous*
- Can model response using Bernoulli dist

$Y_i$	Probability
1	$\Pr\{Y_1 = 1\} = \pi_i$
0	$\Pr\{Y_1 = 0\} = 1 - \pi_i$

- $E\{Y_i\} = \pi_i$
- $Var\{Y_i\} = \pi_i(1 - \pi_i)$

# Goal: express $E\{Y\}$ as function of a covariate $X$

- The simple regression is not appropriate

$$E\{Y_i\} = \beta_0 + \beta_1 X_i$$

It violates several assumptions:

(1) Does not enforce the constraint  $0 \leq E\{Y_i\} \leq 1$  is

(2) Non-normal (binary) distribution of  $\varepsilon \mid X$ :

$$\text{When } Y_i = 0 : \varepsilon_i = 0 - \beta_0 - \beta_1 X_i$$

$$\text{When } Y_i = 1 : \varepsilon_i = 1 - \beta_0 - \beta_1 X_i$$

(3) Non-constant variance

$$\begin{aligned} \text{Var}\{Y_i\} &= \pi_i(1 - \pi_i) \\ &= (\beta_0 + \beta_1 X_i)(1 - \beta_0 - \beta_1 X_i) \end{aligned}$$

# Solution: a Generalized Linear Model

- A generalized linear model is

$$\begin{aligned} E\{Y_i\} &= g(\beta_0 + \beta_1 X_i), \text{ or} \\ g^{-1}(E\{Y_i\}) &= \beta_0 + \beta_1 X_i \end{aligned}$$

where  $g$  is a sigmoid function in  $(0,1)$ .

- $g$  is called the *mean response function*
  - $g^{-1}$  is called the *link function*
- 
- A choice of  $g$  produces different models
    - $g(t) = \text{Identity}$   
→ linear regression
    - $g(t) = \Phi(t) = \text{standard Normal CDF}$   
→ probit regression
    - $g(t) = \frac{\exp(t)}{1+\exp(t)} = \text{CDF of the logistic distrib.}$   
→ logistic regression

# Motivation for Probit

## Regression: Latent Variable

- Assume that the binary response is guided by a non-observed continuous variable
- Example: linear model for blood pressure:

$$bp = \beta_0 + \beta_1 \text{age} + \varepsilon$$

Only observe

$$Y = \begin{cases} 1 \text{ (disease),} & \text{if blood pressure} > c \\ 0 \text{ (healthy),} & \text{if blood pressure} \leq c \end{cases}$$

$$\Pr\{Y = 1\}$$

$$= \Pr\{bp > c\} = \Pr\{\beta_0 + \beta_1 \text{age} + \varepsilon > c\}$$

$$= \Pr\{\varepsilon < \beta_0 + \beta_1 \text{age} - c\}$$

$$= \Pr\left\{\frac{\varepsilon}{\sigma} < \frac{\beta_0 - c}{\sigma} + \frac{\beta_1}{\sigma} \text{age}\right\}$$

$$= \Pr\{z < \beta'_0 + \beta'_1 \text{age}\}$$

$$= \Phi(\beta'_0 + \beta'_1 \text{age})$$

# Logistic Response Function and Logistic Regression

- A sigmoidal response function

$$\begin{aligned} E\{Y_i\} &= \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)} \\ &= \frac{1}{1 + \exp(-\beta_0 - \beta_1 X_i))} \end{aligned}$$

- A monotonic increasing/decreasing function
- Explicit functional form
- Restricts  $0 \leq E(Y_i) \leq 1$
- Example of a **nonlinear** model

- **Logit** link function

$$\log \left( \frac{E\{Y_i\}}{1 - E\{Y_i\}} \right) = \log \left( \frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_1 X_i$$

# Probability Distribution of $Y$ in Logistic Regression

- $Y_i$  are independent but not identically distributed Bernoulli random variables

$$Y_i \stackrel{ind}{\sim} \text{Bernoulli}(\pi_i) \text{ where}$$
$$\pi_i = \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)}$$

– note no more error term!

- Probability density of  $Y_i$

$$f(Y_i) = \pi_i^{Y_i} (1 - \pi_i)^{1-Y_i}$$

- Least Squares Estimates are inappropriate
  - use maximum likelihood for parameter estimation



# Simple Logistic Regression:

## Interpretation of $b_1$

- Fitted value for the individual  $i$

$$- \hat{\pi}_i = \frac{e^{b_0 + b_1 X_i}}{1 + e^{b_0 + b_1 X_i}}$$

- Fitted logistic response (i.e. fitted log odds)

$$\begin{aligned} - \log_e (\widehat{\text{odds}}(X_i)) &= \log_e \frac{\hat{\pi}_i(X_i)}{1 - \hat{\pi}_i(X_i)} \\ &= b_0 + b_1 X_i \end{aligned}$$

- $b_1$  is the slope of the fitted logistic response

- $b_1$  is interpreted as log(odds ratio)

$$\begin{aligned} - b_1 &= \log_e \left( \frac{\hat{\pi}_i(X_i + 1)}{1 - \hat{\pi}_i(X_i + 1)} \right) - \log_e \left( \frac{\hat{\pi}_i(X_i)}{1 - \hat{\pi}_i(X_i)} \right) \\ &= \log_e \left( \frac{\widehat{\text{odds}}(X_i + 1)}{\widehat{\text{odds}}(X_i)} \right) \end{aligned}$$

- $\widehat{\text{odds ratio}}(X_i) = \exp(b_1)$

# Estimation by Maximum Likelihood

- $Y_i \stackrel{ind}{\sim} \text{Bernoulli}(\pi_i)$  where

$$\pi_i = \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)}$$

- Log likelihood:  $\log_e(L) =$

$$\begin{aligned} &= \log \left\{ \prod_{i=1}^n \pi_i^{Y_i} (1 - \pi_i)^{1-Y_i} \right\} \\ &= \sum_{i=1}^n Y_i \log(\pi_i) + \sum_{i=1}^n (1 - Y_i) \log(1 - \pi_i) \\ &= \sum_{i=1}^n Y_i \log\left(\frac{\pi_i}{1 - \pi_i}\right) + \sum_{i=1}^n \log(1 - \pi_i) \\ &= \sum_{i=1}^n Y_i(\beta_0 + \beta_1 X_i) - \sum_{i=1}^n \log(1 + \exp(\beta_0 + \beta_1 X_i)) \end{aligned}$$

- MLEs do not have closed forms

# Equivalent specification:

## Binomial distribution

- Change in notation
  - Data:  $(Y_{ij}, n_i, X_i), i = 1, 2, \dots, c$
  - $X_i$  : predictor for observation  $i$
  - $n_i$  : # of Bernoulli trials in observation  $i$
  - $Y'_i := \sum_{j=1}^{n_i} Y_{ij}$
  - Model:

$$Y'_i \stackrel{\text{ind}}{\sim} \text{Binomial}(n_i, \pi_i), \text{ where}$$

$$\pi_i = \frac{\exp(\beta_0 + \beta_1 X_i)}{1 + \exp(\beta_0 + \beta_1 X_i)}$$

- Log-Likelihood:  $\log_e(L) =$

$$= \log \prod_{i=1}^c \left\{ \binom{n_i}{Y'_i} \pi_i^{Y'_i} (1 - \pi_i)^{n_i - Y'_i} \right\}$$

$$= \sum_{i=1}^c \left\{ Y'_i \log(\pi_i) + (n_i - Y'_i) \log(1 - \pi_i) + \log \left( \binom{n_i}{Y'_i} \right) \right\}$$

$$= \sum_{i=1}^c \left\{ Y'_i \log \frac{\pi_i}{1 - \pi_i} + n_i \log(1 - \pi_i) + \log \left( \binom{n_i}{Y'_i} \right) \right\}$$

# Equivalence of Bernoulli and Binomial Models

- Binomial Log-Likelihood equals Bernoulli Log-Likelihood, up to a constant:

$$\begin{aligned}\log_e(L)^{Binomial} &= \\&= \sum_{i=1}^c \{Y'_i \log(\pi_i) + (n_i - Y'_i) \log(1 - \pi_i)\} + constant \\&= \sum_{i=1}^c \left\{ \sum_{j=1}^{n_i} Y_{ij} \log(\pi_i) + (n_i - \sum_{j=1}^{n_i} Y_{ij}) \log(1 - \pi_i) \right\} + constant \\&= \sum_{i=1}^c \sum_{j=1}^{n_i} \left\{ Y_{ij} \log\left(\frac{\pi_i}{1 - \pi_i}\right) + \log(1 - \pi_i) \right\} + constant \\&= \log_e(L)^{Bernoulli} + constant\end{aligned}$$

- Both models lead to same parameter estimates and inferences, but have different deviances.

# Prospective and Retrospective Studies

- Prospective study: fix predictors, observe the outcome
  - Recruit patients with 2 genotypes, compare occurrence of disease.
  - Expensive; large variance for rare diseases
- Retrospective (=case-control) study: fix outcome, observe predictors
  - Recruit patients with and without the disease, compare the genotypes.
  - Cheaper
- Log odds ratio based on logistic regression is the same for both studies.
  - Not true for other link functions

# Prospective Model With Retrospective Data: Formulation

- Assume a prospective model:

$$\pi(X_i) = P(Y_i = 1|X_i) = \frac{e^{\beta_0 + \beta_1 X_i}}{1 + e^{\beta_0 + \beta_1 X_i}}$$

- However, the subjects are selected retrospectively.
- The random variable is  $Z_i = 1_{\{\text{including subject } i\}}$ , conditional on  $Y$ .
  - Denote  $\theta_0 = P(Z_i = 1|Y_i = 0)$  and  $\theta_1 = P(Z_1 = 1|Y_i = 1)$
  - Note that  $\theta_0$  and  $\theta_1$  are indenpent of  $X_i$ , otherwise introduce selection bias.
- Of interest in retrospective study is  $P(Y_i = 1|X_i, Z_i = 1)$

# Prospective Model With Retrospective Data: Estimation of $\log(\text{OR})$

- With retrospective sampling, we model:

$$\begin{aligned}
 &P(Y_i = 1|X_i, Z_i = 1) \\
 &= \frac{P(Y_i = 1, Z_i = 1|X_i)}{P(Z_i = 1|X_i)} \\
 &= \frac{P(Y_i = 1, Z_i = 1|X_i)}{P(Y_i = 0, Z_i = 1|X_i) + P(Y_i = 1, Z_i = 1|X_i)} \\
 &= \frac{\theta_1 \times \pi(X_i)}{\theta_0 \times \{1 - \pi(X_i)\} + \theta_1 \times \pi(X_i)} \\
 &= \frac{\theta_1 e^{\beta_0 + \beta_1 X_i}}{\theta_0 + \theta_1 e^{\beta_0 + \beta_1 X_i}} = \frac{e^{\log(\theta_1/\theta_0) + \beta_0 + \beta_1 X_i}}{1 + e^{\log(\theta_1/\theta_0) + \beta_0 + \beta_1 X_i}}
 \end{aligned}$$

- Uncover the same  $\beta_1$  as in the prospective study

$$\Rightarrow \log \left\{ \frac{P(Y_i = 1|X_i, Z_i)}{P(Y_i = 0|X_i, Z_i)} \right\} = [\log(\theta_1/\theta_0) + \beta_0] + \beta_1 X_i$$

# Multiple Logistic Regression

- Extension of the response function:

$$E\{Y_i\} = \frac{\exp(\mathbf{X}'_i\beta)}{1 + \exp(\mathbf{X}'_i\beta)}$$

- Extension of the logit link function:

$$\log_e \left( \frac{E\{Y_i\}}{1 - E\{Y_i\}} \right) = \log_e \left( \frac{\pi_i}{1 - \pi_i} \right) = \mathbf{X}'_i\beta$$

- The multivariate likelihood function:

$$\log_e L(\beta) = \sum_{i=1}^n Y_i(\mathbf{X}'_i\beta) - \sum_{i=1}^n \log_e[1 + \exp(\mathbf{X}'_i\beta)]$$

- Interpretation of  $b_j$ :

$$\log_e \left( \frac{\text{odds}(X_j + 1)}{\text{odds}(X_j)} \right)$$

while other predictors are held fixed



# Testing

# Asymptotic Properties of $\hat{\beta}$

- Asymptotic existence and uniqueness:

- $P\{\hat{\beta} \text{ exists and is unique} \rightarrow 1\}$  as  $n \rightarrow \infty$

- Consistency:

- $\hat{\beta} \rightarrow \beta$  as  $n \rightarrow \infty$

- Asymptotic Normality:

- $\hat{\beta} \overset{Ass.}{\sim} \mathcal{N}(\beta, I(\hat{\beta})^{-1})$  as  $n \rightarrow \infty$

- Asymptotic efficiency:

- The MLE has asymptotically smaller variance than many other estimators

# Inference About Individual

## $\beta_j$ : Wald Test

- Test  $H_0 : \beta_j = 0$  versus  $H_a : \beta_j \neq 0$ .

- Test statistic  $z^* = \frac{b_j - 0}{s\{b_j\}}$

- **Approximate** variance  $s^2\{\mathbf{b}\}$

$$s^2\{\mathbf{b}\} = \left( \left[ - \frac{\partial^2 \log_e L(\beta)}{\partial \beta_j \partial \beta_{j'}} \right]_{\beta=\mathbf{b}} \right)^{-1}$$

- **Approximate** distribution of  $z$

–  $z^* \sim \mathcal{N}(0, 1)$ . Alternatively,  $(z^*)^2 \sim \chi_1^2$

– reject  $H_0$  if  $|z^*| > z^{1-\alpha/2}$

– CI for  $\beta_j$ :  $b_j \pm z^{1-\alpha/2} s\{b_j\}$

# Simultaneous Inference

## About Several $\beta_j = 0$ :

### Likelihood Ratio Test

- Multivariate logistic regression

$$\log \left( \frac{E\{Y_i\}}{1 - E\{Y_i\}} \right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1}$$

- Test  $H_0 : \beta_1 = \beta_2 = \cdots = \beta_q = 0$   
versus  $H_a$ : not all  $\beta_1, \beta_2, \cdots, \beta_q = 0$

- Test statistic

$$\begin{aligned} G^2 &= -2 \log_e \left[ \frac{L(\text{reduced model})}{L(\text{full model})} \right] \\ &= -2 [\log_e L(\text{reduced model}) - \log_e L(\text{full model})] \end{aligned}$$

- **Approximate** distr of  $G^2$  for large  $n$ 
  - Reject  $H_0$  if  $G^2 > \chi^2(1 - \alpha, q)$

# Comments: Wald Test Versus Likelihood Ratio Test

- Both tests are approximate, for large  $n$
- Likelihood Ratio test:  $\beta_j = 0$ , or several  $\beta_j = 0$  simultaneously
  - no other values of  $\beta_j$
  - no one-sided tests
- Wald test:  $\beta_j = \beta_j^{of\ interest}$ , or linear combinations of  $\beta_j$
- When testing a single  $H_0 : \beta_j = 0$ , the tests may lead to different conclusions
  - due to the approximate nature of the tests
  - unlike in linear regression

# **Quality of Fit: Replicated data**

# Lack of Fit for Replicated Data: Pearson $\chi^2$

$$H_0 : \log \left( \frac{E\{Y_{ij}\}}{1 - E\{Y_{ij}\}} \right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1} \text{ vs}$$

$$H_a : \log \left( \frac{E\{Y_{ij}\}}{1 - E\{Y_{ij}\}} \right) \neq \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1}$$

- $c$  covariate configurations, each with  $n_j$  cases
- Observed counts
  - $O_{i0}$ : observed # of 0's in configuration  $i$
  - $O_{i1}$ : observed # of 1's in configuration  $i$
- Expected counts
  - $E_{i0} = n_i(1 - \hat{\pi}_i)$ : expected # of 0's in conf.  $i$
  - $E_{i1} = n_i(\hat{\pi}_i)$ : expected # of 1's in conf.  $i$
- Test statistic: reject  $H_0$  if

$$X^2 = \sum_{i=1}^c \sum_{k=0}^1 \frac{(O_{jk} - E_{jk})^2}{E_{jk}} > \chi^2(1 - \alpha, c - p)$$

# Lack of Fit for Replicated Data: Deviance

$$H_0 : \quad \log \left( \frac{E\{Y_{ij}\}}{1 - E\{Y_{ij}\}} \right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1} \quad \text{vs}$$

$$H_a : \quad \log \left( \frac{E\{Y_{ij}\}}{1 - E\{Y_{ij}\}} \right) \neq \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1}$$

- $c$  covariate configurations, each with  $n_i$  cases

- $H_0$ : model of interest;  $H_a$ : saturated model

- Model of interest

- $E\{Y_{ij}\} = \pi_i$ ;  $E\{\widehat{Y_{ij}}\} = \hat{\pi}_i$

- Saturated model

- $E\{Y_{ij}\} = p_i$ ;  $E\{\widehat{Y_{ij}}\} = \hat{p}_i = \frac{\sum_j Y_{ij}}{n_i}$

- Test statistic: LR, also called deviance

- Deviance of a saturated model always = 0



# Lack of Fit for Replicated Data: Deviance

$$\begin{aligned}
 G^2 &= DEV(X_0, X_1, \dots, X_{p-1}) \\
 &= -2 [ \log_e L(\text{current model}) - \log_e L(\text{saturated model}) ] \\
 &= -2 \sum_{i=1}^c \left[ \sum_{j=1}^{n_i} Y_{ij} \log_e \hat{\pi}_i + (n_i - \sum_{j=1}^{n_i} Y_{ij}) \log_e (1 - \hat{\pi}_i) \right] \\
 &\quad + 2 \sum_{i=1}^c \left[ \sum_{j=1}^{n_i} Y_{ij} \log_e \hat{p}_i + (n_i - \sum_{j=1}^{n_i} Y_{ij}) \log_e (1 - \hat{p}_i) \right] \\
 &= -2 \sum_{i=1}^c \left[ \sum_{j=1}^{n_i} Y_{ij} \log_e \left( \frac{\hat{\pi}_i}{\hat{p}_i} \right) + (n_i - \sum_{j=1}^{n_i} Y_{ij}) \log_e \left( \frac{1 - \hat{\pi}_i}{1 - \hat{p}_i} \right) \right]
 \end{aligned}$$

- Reject  $H_0$  if  $G^2 > \chi^2(1 - \alpha, c - p)$
- Approximation  $\chi^2(1 - \alpha, c - p)$  can be poor
  - The closer the distribution to Gaussian, and the closer the link to identity, the better the approximation
  - Unlike with the LR test, the quality of approximation does not improve with the sample size

## Note: Can Use Deviance for LR Test of Nested Models

$$\log \left( \frac{E\{Y_{ij}\}}{1 - E\{Y_{ij}\}} \right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1}$$

- Test  $H_0 : \beta_1 = \beta_2 = \cdots = \beta_q = 0$   
versus  $H_a$ : not all  $\beta_1, \beta_2, \cdots, \beta_q = 0$
- Likelihood Ratio test statistic  $G^2 =$ 
$$\begin{aligned} &= -2 \log_e \left[ \frac{L(\text{reduced model})}{L(\text{full model})} \right] \\ &= -2 [\log_e L(\text{reduced model}) - \log_e L(\text{full model})] \\ &= -2 [\log_e L(\text{reduced model}) - \log_e L(\text{saturated model})] \\ &\quad + 2 [\log_e L(\text{full model}) - \log_e L(\text{saturated model})] \\ &= \text{Deviance}(\text{reduced model}) - \text{Deviance}(\text{full model}) \end{aligned}$$
- Approximate distr of  $G^2$  for large  $n$
- Reject  $H_0$  if  $G^2 > \chi^2(1 - \alpha, q)$

# **Quality of Fit: Individual Observations**

# Non-Replicated Data:

## Hosmer-Lemeshow

### Goodness of Fit

$$H_0 : \log \left( \frac{E\{Y_i\}}{1 - E\{Y_i\}} \right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1} \text{ vs}$$
$$H_a : \log \left( \frac{E\{Y_i\}}{1 - E\{Y_i\}} \right) \neq \beta_0 + \beta_1 X_{i1} + \cdots + \beta_{p-1} X_{i,p-1}$$

- Group cases based on values of estimated probabilities  $\hat{\pi}_i$  into  $c$  groups
  - E.g., find  $c = 9$  groups based on percentiles
- Apply Pearson  $\chi^2$  test to the groups
- Reject  $H_0$  if  $X^2 > \chi^2(1 - \alpha, c - 2)$ 
  - showed by simulation that this distribution is appropriate

# Can Write Pearson $\chi^2$ (But Not Use for Tests)

- Test statistic:

$$\begin{aligned} X^2 &= \sum_{i=1}^c \sum_{k=0}^1 \frac{(O_{ik} - E_{ik})^2}{E_{ik}} \\ &= \sum_{i=1}^c \frac{(O_{i0} - E_{i0})^2}{E_{i0}} + \sum_{i=1}^c \frac{(O_{i1} - E_{i1})^2}{E_{i1}} \end{aligned}$$

- The corresponding quantities (KNNL p. 591):

- $c = n, n_i = 1$
- $O_{i0} = 1 - Y_i, O_{i1} = Y_i$
- $E_{i0} = 1 - \hat{\pi}_i, E_{i1} = \hat{\pi}_i$

- Test statistic:

$$\begin{aligned} X^2 &= \sum_{i=1}^n \frac{[(1 - Y_i) - (1 - \hat{\pi}_i)]^2}{1 - \hat{\pi}_i} + \sum_{i=1}^n \frac{(Y_i - \hat{\pi}_i)^2}{\hat{\pi}_i} \\ &= \sum_{i=1}^n \frac{(Y_i - \hat{\pi}_i)^2}{1 - \hat{\pi}_i} + \sum_{i=1}^n \frac{(Y_i - \hat{\pi}_i)^2}{\hat{\pi}_i} = \sum_{i=1}^n \frac{(Y_i - \hat{\pi}_i)^2}{\hat{\pi}_i(1 - \hat{\pi}_i)} \end{aligned}$$

# Can Write Deviance (But Not Use for Tests)

- Test statistic  $G^2 = DEV(X_0, X_1, \dots, X_{p-1})$ :

$$= -2 \sum_{i=1}^c \left[ \sum_{j=1}^{n_i} Y_{ij} \log \left( \frac{\hat{\pi}_i}{\hat{p}_i} \right) + (n_i - \sum_{j=1}^{n_i} Y_{ij}) \log \left( \frac{1 - \hat{\pi}_i}{1 - \hat{p}_i} \right) \right]$$

- The corresponding quantities (KNNL p. 592):

$$- c = n, \quad n_i = 1, \quad \sum_{j=1}^{n_i} Y_{ij} = Y_i, \quad \hat{p}_i = \sum_{j=1}^{n_i} Y_{ij} / n_i = Y_i$$

- Test statistic:

$$\begin{aligned} G^2 &= -2 \sum_{i=1}^n \left[ Y_i \log \left( \frac{\hat{\pi}_i}{Y_i} \right) + (1 - Y_i) \log \left( \frac{1 - \hat{\pi}_i}{1 - Y_i} \right) \right] \\ &= -2 \sum_{i=1}^n [ Y_i \log \hat{\pi}_i + (1 - Y_i) \log(1 - \hat{\pi}_i) \\ &\quad - Y_i \log Y_i - (1 - Y_i) \log(1 - Y_i) ] \\ &= -2 \sum_{i=1}^n [ Y_i \log \hat{\pi}_i + (1 - Y_i) \log(1 - \hat{\pi}_i) ] \end{aligned}$$

# Diagnostics: Residuals

- Logistic Regression Residuals

$$e_i = \begin{cases} 1 - \hat{\pi}_i, & \text{if } Y_i = 1 \\ -\hat{\pi}_i, & \text{if } Y_i = 0 \end{cases}$$

- Pearson residual

$$r_{P_i} = \frac{Y_i - \hat{\pi}_i}{\sqrt{\hat{\pi}_i(1 - \hat{\pi}_i)}}$$

- $e_i$  divided by the standard error of  $Y_i$
- $\sum_{i=1}^n r_{P_i}^2$  equals **non-replicated** Pearson  $X^2$

- Studentized Pearson Residual

$$r_{P_i} = \frac{Y_i - \hat{\pi}_i}{\sqrt{\hat{\pi}_i(1 - \hat{\pi}_i) \cdot (1 - h_{ii})}}$$

- $e_i$  divided by the SE of  $e_i \rightarrow$  unit variance
- $h_{ii}$  is the diagonal element of the hat matrix

$$\mathbf{H} = \hat{\mathbf{W}}^{\frac{1}{2}} \mathbf{X} (\mathbf{X}' \hat{\mathbf{W}} \mathbf{X})^{-1} \mathbf{X}' \hat{\mathbf{W}}^{\frac{1}{2}}, \text{ where}$$

$$\hat{\mathbf{W}} = \text{diag}(\hat{\pi}_i(1 - \hat{\pi}_i))$$

# Diagnostics: Residuals

- Deviance Residuals

$$\begin{aligned}d_i &= \text{sign}(Y_i - \hat{\pi}_i) \sqrt{-2 \left[ Y_i \log \frac{\hat{\pi}_i}{Y_i} + (1 - Y_i) \log \frac{1 - \hat{\pi}_i}{1 - Y_i} \right]} \\&= \text{sign}(Y_i - \hat{\pi}_i) \sqrt{-2 [Y_i \log \hat{\pi}_i + (1 - Y_i) \log(1 - \hat{\pi}_i)]}\end{aligned}$$

- the signed square root of the contribution of  $Y_i$  to the model deviance

- Analysis of residuals

- unknown distribution of residuals under true model
- plot residual by predicted value. A flat lowess smooth to this plot suggests good model

- Other summaries as in linear regression

- DFFITS, DFBETAS
- $\Delta\chi^2$ ,  $\Delta\text{dev}$ , Cook's distance (see KNNL p. 598)



# Graphical Check of the Fit

- Partition the observations into groups of covariate patterns  $x_i$ .

- Haldane (1956) recommended to plot

$$\hat{\eta}_i = \log \frac{y_i + 0.5}{n_i - y_i + 0.5}$$

against covariate patterns  $\mathbf{x}_i$ .

- The plot should be roughly linear if the model is appropriate for the data
- When all  $n_i = 1$  or all  $n_i$  are small, one can group the data with nearby  $x$  values to make the plot

# Overdispersion

# Overdispersion

$$Y \stackrel{ind}{\sim} \text{Binomial}(n, \pi), \quad \pi = \frac{\exp(\mathbf{X}'\beta)}{1 + \exp(\mathbf{X}'\beta)}$$

- Implies  $E\{Y\} = \pi$ ,  $Var\{Y\} = n\pi(1 - \pi)$ 
  - Overdispersion:  $Var\{Y\} > n\pi(1 - \pi)$
  - Underdispersion:  $Var\{Y\} < n\pi(1 - \pi)$
- Mechanisms of overdispersion:
  - Suppose  $Y_1, Y_2, \dots, Y_n$  are Bernoulli r.v.,  $E\{Y_i\} = \pi$ .
  - define  $Y = \sum_{i=1}^n Y_i$
  - Can think of at least two situations when  $Y$  does not have a Binomial distribution (and therefore a different variance)

# Overdispersion from correlation

- Suppose  $Y_1, Y_2, \dots, Y_n$  are not independent
- Suppose all pairs  $(Y_i, Y_j)$  have a same correlation  $\rho$

$$\begin{aligned} \text{Var}(Y) &= \text{Cov} \left( \sum_{i=1}^n Y_i, \sum_{i=1}^n Y_i \right) \\ &= \sum_{i=1}^n \text{Var}(Y_i) + \sum_{i \neq j} \text{Corr}(Y_i, Y_j) \sqrt{\text{Var}(Y_i)} \sqrt{\text{Var}(Y_j)} \\ &= n\pi(1 - \pi) + n(n - 1)\rho\pi(1 - \pi) \\ &> n\pi(1 - \pi) \end{aligned}$$

- The variance exceeds the variance of the Binomial distribution

# Overdispersion from clustered data

- Suppose  $Y = \sum_{i=1}^n Y_i \mid \pi \sim \text{Binomial}(\pi)$
- Suppose  $\pi$  is a random variable,  
 $E\{\pi\} = p$ ,  $\text{Var}\{\pi\} = p(1 - p)$ 
  - Special case:  $\pi \sim \text{Beta}(\alpha, \beta)$

$$\begin{aligned} E\{Y\} &= E\{E\{Y \mid \pi\}\} = p \\ \text{Var}\{Y\} &= \text{Var}\{E\{Y \mid \pi\}\} + E\{\text{Var}\{Y \mid \pi\}\} \\ &= \text{Var}\{n\pi\} + E\{n\pi(1 - \pi)\} \\ &= n^2 \text{Var}\{\pi\} + np - n[\text{Var}\{\pi\} + p^2] \\ &= np(1 - p) + n(n - 1)\text{Var}\{\pi\} \\ &> np(1 - p) \end{aligned}$$

- The variance exceeds the variance of the Binomial distribution

# Modeling Overdispersion

- Introduce additional parameter

$$E\{Y_i\} = n_i\pi_i, \quad Var\{Y_i\} = \phi n_i\pi_i\{1 - \pi_i\}$$

- When  $\phi \neq 1$ ,  $Y_i$  follows a *quasi-binomial distribution*. The distribution is characterized by its expectation and variance. The probability distribution function is unspecified.

- $\hat{\phi} = \chi^2 / (n - p)$

- $\chi^2$  is the Pearson lack of fit statistic (= sum of squared Pearson residuals with **non-replicated** data)
- $p$  is the number of parameters in the model
- $\hat{\phi} \gg 1$  indicates evidence of overdispersion.

- Since  $\phi$  does not affect  $E\{Y_i\}$ , modeling overdispersion does not change  $\hat{\beta}$ .

- $SE\{\hat{\beta}\}$  is multiplied by  $\sqrt{\hat{\phi}}$ .

# Comparing Nested Models in Presence of Overdispersion

- Regular likelihood-based approaches (e.g. LRT, AIC) are not applicable.
- F test approximates deviance-based LR test

$$F = \frac{D_{reduced} - D_{full}}{df_{reduced} - df_{full}} / \hat{\phi} \stackrel{ass.}{\sim} H_0 F_{df_{reduced}-df_{full}, df_{full}}$$

- Assumes roughly equal covariate classes
- Modeling strategy:
  - fit the full model (with all predictors)
  - estimate  $\hat{\phi}$
  - compare nested models with  $F$  test to reduce the number of predictors

# **Prediction and Classification**



# Prediction of the Mean

- Point estimate for the link  $\hat{\pi}'_h$ :

$$\hat{\pi}'_h = \mathbf{X}'_h \mathbf{b}$$

- Point estimate for the response  $\hat{\pi}_h$ :

$$\hat{\pi}_h = \frac{1}{1 + \exp(-\hat{\pi}'_h)} = \frac{1}{1 + \exp(-\mathbf{X}'_h \mathbf{b})}$$

- Interval estimate for the link  $\hat{\pi}'_h$ :

$$s^2\{\hat{\pi}'_h\} = \mathbf{X}'_h s^2\{\mathbf{b}\} \mathbf{X}_h$$

$$(1 - \alpha)\% \text{ CI for } \hat{\pi}'_h : \hat{\pi}'_h \pm z^{1-\alpha/2} s\{\hat{\pi}'_h\} = (L, U)$$

- Approximate interval estimate for the response  $\hat{\pi}_h$ :

$$(1 - \alpha)\% \text{ CI for } \hat{\pi}_h : \left( \frac{1}{1 + \exp(-L)}, \frac{1}{1 + \exp(-U)} \right)$$

- Use Bonferroni if multiple  $\mathbf{X}$  are of interest

# Measures of Agreement

- Have  $N$  observations
- Consider all pairs of distinct responses
  - In this example  $t = 16 \times 14 = 224$
- Compare predicted probabilities
  - Concordant if  $\hat{\pi}_{Y=1} > \hat{\pi}_{Y=0}$
  - Discordant if  $\hat{\pi}_{Y=1} < \hat{\pi}_{Y=0}$
  - Tie if  $\hat{\pi}_{Y=1} = \hat{\pi}_{Y=0}$
- Measures of agreement
  - Somers' D :  $(\#C - \#D)/t$
  - Goodman-Kruskal Gamma :  $(\#C - \#D)/(\#C + \#D)$
  - Kendall's Tau-a :  $(\#C - \#D)/(.5N(N-1))$
  - c :  $(\#C + .5(t - \#C - \#D))/t$

# Prediction of a New Observation (i.e. Classification)

- Choose a cutoff  $c \in (0, 1)$

$$\hat{Y}_h = \begin{cases} 1, & \text{if } \hat{\pi}_h > c \\ 0, & \text{if } \hat{\pi}_h \leq c \end{cases}$$

- Sensitivity:

$$\frac{\# \text{ predicted '1' \& true '1'}}{\# \text{ true '1'}} = \frac{\sum_{i=1}^n \hat{Y}_i \cdot Y_i}{\sum_{i=1}^n Y_i}$$

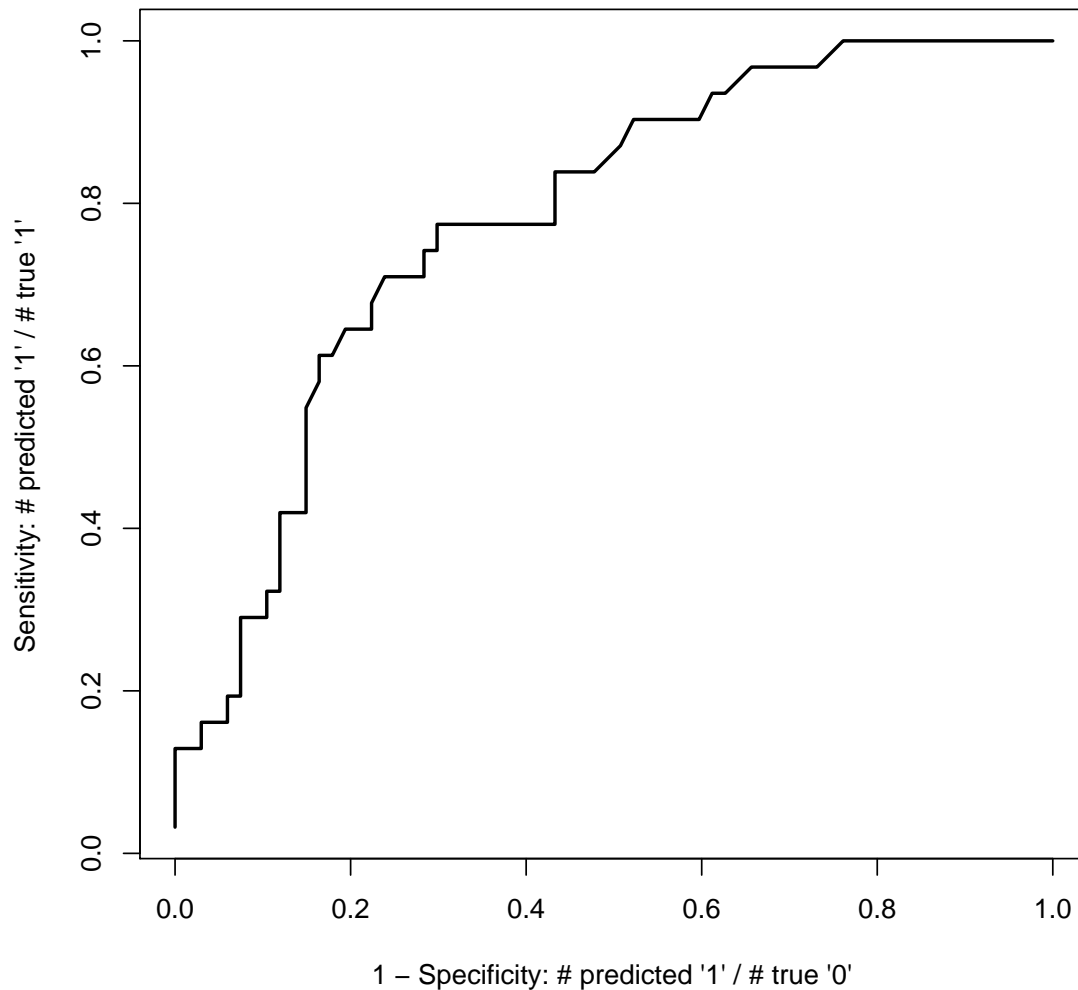
- Specificity:

$$\frac{\# \text{ predicted '0' \& true '0'}}{\# \text{ true '0'}} = \frac{\sum_{i=1}^n (1 - \hat{Y}_i) \cdot (1 - Y_i)}{\sum_{i=1}^n (1 - Y_i)}$$

- Vary the cut-off  $c \in (0, 1)$ , and choose  $c$  to optimize sensitivity and specificity

# ROC curve for classification

Vary  $c$ , and plot sensitivity vs 1-specificity.  
Evaluate models by area under the curve.



# Evaluation of the Predictive Ability of the Model

- Area under ROC can be used to compare models
  - Area = 1  $\rightarrow$  perfect classification
  - Area = .5  $\rightarrow$  random classification
- Classification on the training set is overly optimistic
- Use cross-validation to construct a more accurate ROC curve

# **Variable Selection**

# Automatic Variable Selection

- Exhaustive search. Minimize:

$$-2 \log_e L(\mathbf{b})$$

$$AIC_p = -2 \log_e L(\mathbf{b}) + 2p$$

$$BIC_p = -2 \log_e L(\mathbf{b}) + p \log_e(n)$$

- Heuristic search
  - forward selection; backward elimination; step-wise selection
  - based on Wald statistic and Normal distribution

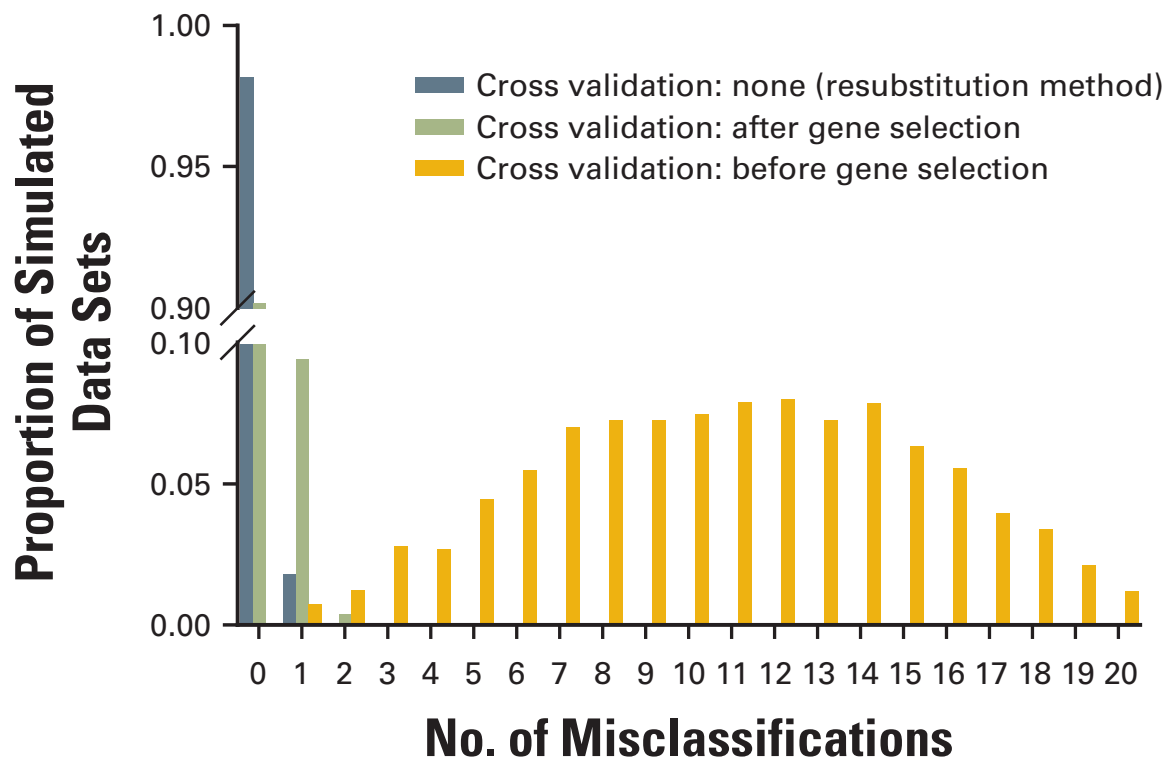
# Variable Selection Should be Done as Part of Cross-Validation

- Example from Simon *et al.*, JNCI, 2003.
- Simulated data with no structure
  - 20 observations with random labels
  - 6,000 possible but unrelated predictors
  - Repeated 200 times
- Estimated predictive accuracy using
  - no cross-validation
  - selecting features on full dataset, then using cross-validation
  - selecting features at each step of cross-validation



# Variable Selection Should be Done as Part of Cross-Validation

Example from Simon *et al.*, JNCI, 2003.



- Conclusion

- Incorporating selection of predictors within the cross-validation procedure is key