

# Likelihood ratio tests

• Department Seminar: Dr. Mine Getinkaya-Rundel  
September 26, 12pm - 1pm Hirby 120

- Sign up to meet with  
the speaker (11 - 11:30)
- Solutions posted to Friday's class activity

## Last time

Data on the RMS *Titanic* disaster. We have data on 891 passengers on the ship, with the following variables:

- + Passenger: A unique ID number for each passenger.
- + Survived: An indicator for whether the passenger survived (1) or perished (0) during the disaster.
- + Pclass: Indicator for the class of the ticket held by this passengers; 1 = 1st class, 2 = 2nd class, 3 = 3rd class.
- + Sex: Binary Indicator for the biological sex of the passenger.
- + Age: Age of the passenger in years; Age is fractional if the passenger was less than 1 year old.
- + Fare: How much the ticket cost in US dollars.
- + others

## Last time

$$Sex = \begin{cases} 1 & \text{male} \\ 0 & \text{female} \end{cases}$$

*Is there a relationship between passenger age and their probability of survival, after accounting for sex, passenger class, and the cost of their ticket?*

$\beta_1$  = slope for female passengers

$\beta_1 + \beta_3$  = slope for males

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 Age_i + \beta_2 Sex_i + \beta_3 Age_i \cdot Sex_i + \beta_4 \log(Fare_i + 1)$$

What hypotheses should we test to investigate this research question?

$$H_0: \beta_1 = \beta_3 = 0$$

$$H_A: \text{at least one of } \beta_1, \beta_3 \neq 0$$

## wald test

1) Rewrite the order of the terms

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{Sex}_i + \beta_2 \log(\text{Fare}_i + 1) + \beta_3 \text{Age}_i + \beta_4 \text{Age}_i \cdot \text{Sex}_i$$

2)  $\beta = \begin{bmatrix} \beta_{(1)} \\ \beta_{(2)} \end{bmatrix} \quad \beta_{(2)} = \begin{bmatrix} \beta_3 \\ \beta_4 \end{bmatrix}$

$H_0: \beta_{(2)} = \beta_{(2)}^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$H_A: \beta_{(2)} \neq \beta_{(2)}^0$

3) Partition variance matrix :  $\Sigma^{-1}(\beta) = \begin{bmatrix} \Sigma^{11} & \Sigma^{12} \\ \Sigma^{21} & \Sigma^{22} \end{bmatrix}$

 $\text{Var}(\hat{\beta}_{(2)}) = \Sigma^{22}$

$$\Sigma^{-1}(\beta) = (k+1) \times (k+1) \text{ matrix}$$

Here:  $5 \times 5$  matrix

$$\Rightarrow \Sigma^{11} \in \mathbb{R}^{3 \times 3}$$

$$\Sigma^{12} \in \mathbb{R}^{3 \times 2}$$

$$\Sigma^{21} \in \mathbb{R}^{2 \times 3}$$

$$\Sigma^{22} \in \mathbb{R}^{2 \times 2}$$

$$\Sigma^{22} \in \mathbb{R}^{2 \times 2} \quad (\text{b/c } \beta_{(2)} \in \mathbb{R}^2)$$

$$4) \text{ Test Statistic: } W = (\hat{\beta}_{(2)} - \beta_{(2)}^0)^T (Z^{22})^{-1} (\hat{\beta}_{(2)} - \beta_{(2)}^0)$$

5) p-value, under  $H_0$ ,  $W \sim \chi^2_q$  # parameters tested

$$\text{Hence: } W \sim \chi^2_2$$

$$\underbrace{\log\left(\frac{p_i}{1-p_i}\right)}_{\text{Full model}} = \beta_0 + \beta_1 \text{Sex}_i + \beta_2 \log(\text{Fare}_i + 1) + \beta_3 \text{Age}_i + \beta_4 \text{Age}_i \cdot \text{Sex}_i$$

$$H_0: \beta_3 = \beta_4 = 0$$

$$\Rightarrow \underbrace{\text{Reduced model}}_{\text{ }}: \log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{Sex}_i + \beta_2 \log(\text{Fare}_i + 1)$$

# Likelihood ratio tests

```
...  
## Coefficients:  
##                                     Estimate Std. Error z value Pr(>|z|)  
## (Intercept)      -1.40695     0.44682 -3.149   0.00164 **  
## Age             0.01107     0.01107  1.000   0.31730  
## Sexmale        -1.27467     0.41654 -3.060   0.00221 **  
## log(Fare + 1)   0.69449     0.11065  6.276 3.47e-10 ***  
## Age:Sexmale    -0.03638     0.01378 -2.639   0.00831 **  
##  
## Null deviance: 964.52 on 713 degrees of freedom  
## Residual deviance: 697.21 on 709 degrees of freedom  
...  
Deviance!
```

What information replaces  $R^2$  and  $R^2_{adj}$  in the GLM output?

For logistic:  $AIC = 2(k+1) + \text{deviance}$

## Deviance

(residual deviance)

Deviance similar to SSE  
for linear regression

**Definition:** The *deviance* of a fitted model with parameter estimates  $\hat{\beta}$  is given by

(want to minimize deviance)

$$2\ell(\text{saturated model}) - 2\ell(\hat{\beta})$$

still  $\sum_{i=1}^n \log(\hat{p}_i^{y_i} (1-\hat{p}_i)^{1-y_i})$

But we estimate  $\hat{p}_i$  separately  
for each point

Saturated model:  $\hat{p}_i = y_i$

$$\Rightarrow \sum_{i=1}^n \log(1) = 0$$

For binary logistic regression: deviance =  $-2\ell(\hat{\beta})$

log likelihood for the fitted model

$$\sum_{i=1}^n \log(\hat{p}_i^{y_i} (1-\hat{p}_i)^{1-y_i})$$

e.g.

$$\hat{p}_i = \frac{e^{\hat{\beta}^T x_i}}{1 + e^{\hat{\beta}^T x_i}}$$

# Residual and null deviance

```
m1 <- glm(Survived ~ Age*Sex + log(Fare + 1),  
           data = titanic, family = binomial)  
summary(m1)
```

...  
## Null deviance: 964.52 on 713 degrees of freedom  
## Residual deviance: 697.21 on 709 degrees of freedom

...

Residual deviance: deviance for

$$\Rightarrow -2 \ell(\hat{\beta}) = 697.21$$

$$\Rightarrow \ell(\hat{\beta}) = -\frac{697.21}{2}$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{Sex}_i + \beta_2 \log(\text{Fare}_i + 1) + \beta_3 \text{Age}_i + \beta_4 \text{Age}_i \cdot \text{Sex}_i$$

Null deviance: deviance for  
(intercept-only model)

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 \quad (\Rightarrow \frac{e^{\beta_0}}{1+e^{\beta_0}} = \bar{Y} = \text{prevalence of } (s))$$

Linear regression : (Intuition)

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

SSTotal



null deviance

SSE



residual  
deviance

+ SSReg  
(SSModel)

Drop-in-deviance:

## Comparing deviances

$$G = 708.04 - 697.21 \\ = 10.83$$

```
m1 <- glm(Survived ~ Age*Sex + log(Fare + 1),  
           data = titanic, family = binomial)  
summary(m1)
```

Full model

```
...  
## Null deviance: 964.52 on 713 degrees of freedom  
## Residual deviance: 697.21 on 709 degrees of freedom
```

Reduced model

```
m2 <- glm(Survived ~ Sex + log(Fare + 1),  
           data = titanic, family = binomial)  
summary(m2)
```

...

```
## Null deviance: 964.52 on 713 degrees of freedom  
## Residual deviance: 708.04 on 711 degrees of freedom
```

...

# Comparing deviances

Full model:

Hypotheses:

Reduced model:

Test statistic:

# Comparing deviances

Full model:

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 Sex_i + \beta_2 \log(Fare_i + 1) + \beta_3 Age_i + \beta_4 Age_i \cdot Sex_i$$

Reduced model:

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 Sex_i + \beta_2 \log(Fare_i + 1)$$

$$G = 2\ell(\hat{\beta}) - 2\ell(\hat{\beta}^0)$$

Why is  $G$  always  $\geq 0$ ?

# Comparing deviances

Full model:

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 Sex_i + \beta_2 \log(Fare_i + 1) + \beta_3 Age_i + \beta_4 Age_i \cdot Sex_i$$

Reduced model:

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 Sex_i + \beta_2 \log(Fare_i + 1)$$

$$G = 2\ell(\hat{\beta}) - 2\ell(\hat{\beta}^0) = 10.83$$

If the reduced model is correct, how unusual is  $G = 10.83$ ?

# Likelihood ratio test