

Lecture 2: Fitting and interpreting logistic regression models

Ciaran Evans

Last time: Dengue data

Data: Data on 5720 Vietnamese children, admitted to the hospital with possible dengue fever. Variables include:

- ▶ *Sex*: patient's sex (female or male)
- ▶ *Age*: patient's age (in years)
- ▶ *WBC*: white blood cell count
- ▶ *PLT*: platelet count
- ▶ other diagnostic variables. . .
- ▶ *Dengue*: whether the patient has dengue (0 = no, 1 = yes)

Logistic regression model

$$p_i = \mathbb{E}[Y_i | X_i]$$

$$Y_i \sim \text{Bernoulli}(p_i)$$

(random component)

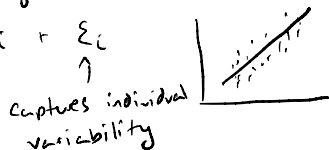
$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{WBC}_i \quad (\text{systematic component})$$

Why is there no noise term ε_i in the logistic regression model?

Discuss for 1-2 minutes with your neighbor, then we will discuss as a class.

Two options for writing linear regression model: $\beta_0 + \beta_1 X$

option 1: $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$



option 2: If we assume $\varepsilon_i \sim N(0, \sigma^2)$, then

$$\left. \begin{array}{l} \text{no } \varepsilon_i \text{ here} \\ \text{b/c captures randomness} \end{array} \right\} \begin{array}{l} Y_i \sim N(\mu_i, \sigma^2) \quad (\text{random}) \\ \mu_i = \beta_0 + \beta_1 X_i \quad (\text{systematic}) \end{array}$$

Fitting the logistic regression model

$$Y_i \sim \text{Bernoulli}(p_i)$$

"generalized
linear model"



$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{WBC}_i$$

```
m1 <- glm(Dengue ~ WBC, data = dengue,  
           family = binomial)  
summary(m1)
```

formula (in R)

family specifies distribution of response variable
logistic regression: family = binomial
linear regression: family = gaussian
poisson " " " = poisson

Fitting the logistic regression model

$$Y_i \sim \text{Bernoulli}(p_i)$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{WBC}_i$$

```
m1 <- glm(Dengue ~ WBC, data = dengue,  
           family = binomial)  
summary(m1)
```

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = 1.737 - 0.361 \text{WBC}_i$$

```
##
```

```
## Call:
```

```
## glm(formula = Dengue ~ WBC, family = binomial, data = de
```

```
##
```

```
## Coefficients:
```

	Estimate	Std. Error	z value	Pr(> z)
## (Intercept)	1.73743	0.08499	20.44	<2e-16 ***
## WBC	-0.36085	0.01243	-29.03	<2e-16 ***

```
##
```

z-values instead of
t-values

Making predictions

$$\log \left(\frac{\hat{p}_i}{1 - \hat{p}_i} \right) = 1.737 - 0.361 \text{ WBC}_i$$

Work in groups of 2-3 on the following questions:

- ▶ What is the predicted odds of dengue for a patient with a WBC of 10?
- ▶ For a patient with a WBC of 10, is the predicted probability of dengue > 0.5 , < 0.5 , or $= 0.5$?
- ▶ What is the predicted *probability* of dengue for a patient with a WBC of 10?

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = 1.737 - 0.361 WBC_i$$

$$WBC = 10$$

$$\log odds = 1.737 - 0.361(10) = -1.873$$

$$odds = e^{-1.873} = \underline{0.154}$$

$$\text{If } p = 0.5 \quad \Leftrightarrow \quad odds = 1 \quad \Leftrightarrow \quad \log odds = 0$$

$$p < 0.5 \quad \Leftrightarrow \quad odds < 1$$

$$p > 0.5 \quad \Leftrightarrow \quad odds > 1$$

\Rightarrow probability < 0.5

$$\hat{p} = \frac{e^{-1.873}}{1 + e^{-1.873}} \approx 0.133$$

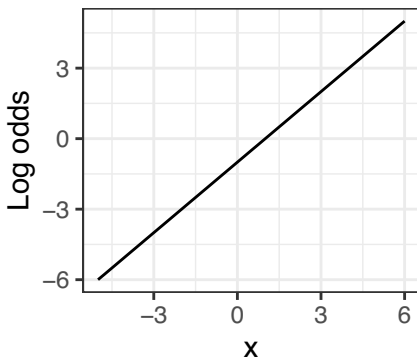
$$\frac{p}{1-p} = odds$$

$$p = \frac{odds}{1 + odds}$$

$$p = \frac{e^{B_0 + B_1 X}}{1 + e^{B_0 + B_1 X}} = \frac{e^{\log odds}}{1 + e^{\log odds}}$$

Shape of the regression curve

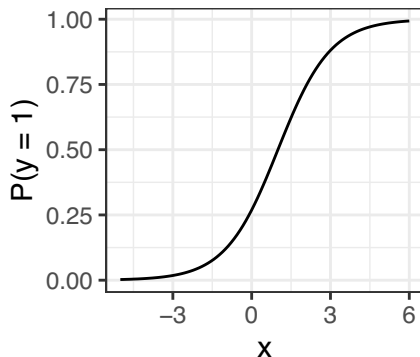
$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_i$$



($\beta_1 > 0$)

\swarrow

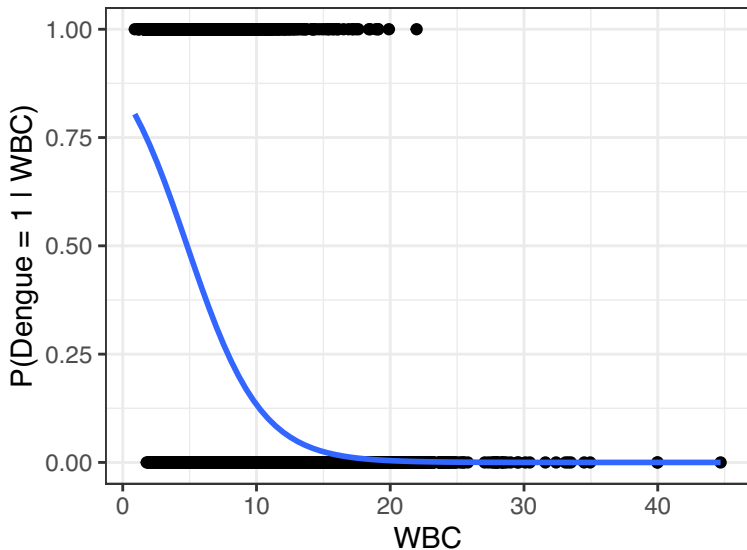
$$p_i = \frac{e^{\beta_0 + \beta_1 X_i}}{1 + e^{\beta_0 + \beta_1 X_i}}$$



$$= p(Y_i = 1)$$

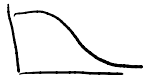
$$Y_i \sim \text{Bernali}(p)$$

Plotting the fitted model for dengue data



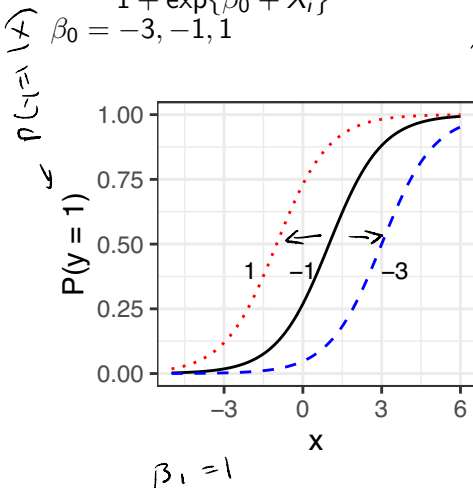
Shape of the regression curve

$$\beta_1 = -1$$

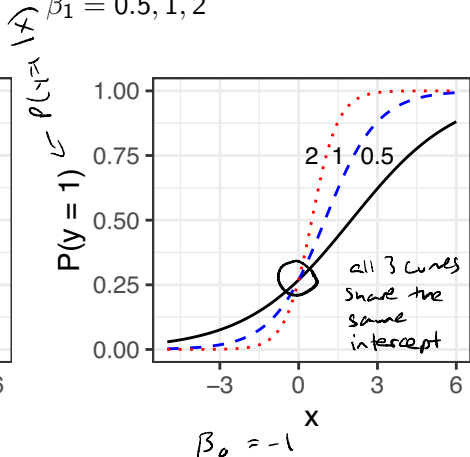


How does the shape of the fitted logistic regression depend on β_0 and β_1 ?

$$p_i = \frac{\exp\{\beta_0 + X_i\}}{1 + \exp\{\beta_0 + X_i\}} \text{ for } \beta_0 = -3, -1, 1$$



$$p_i = \frac{\exp\{-1 + \beta_1 X_i\}}{1 + \exp\{-1 + \beta_1 X_i\}} \text{ for } \beta_1 = 0.5, 1, 2$$



Interpretation

$$\log \left(\frac{\hat{p}_i}{1 - \hat{p}_i} \right) = 1.737 - 0.361 \text{ WBC}_i$$

Work in groups of 2-3 for on the following questions:

- ▶ Are patients with a higher WBC more or less likely to have dengue?
- ▶ What is the change in *log odds* associated with a unit increase in WBC?
- ▶ What is the change in *odds* associated with a unit increase in WBC?

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = 1.737 - 0.361 \text{ WBC}_i$$

$$\text{WBC} \rightarrow \text{WBC} + 1$$

$$\begin{aligned} \log \text{odds} (\text{WBC} + 1) &= 1.737 - 0.361 (\text{WBC} + 1) \\ - \log \text{odds} (\text{WBC}) &= 1.737 - 0.361 (\text{WBC}) \\ &= -0.361 \quad (\log \text{odds decreases by } 0.361) \end{aligned}$$

$$\begin{aligned} \frac{\text{odds} (\text{WBC} + 1)}{\text{odds} (\text{WBC})} &= \frac{e^{1.737 - 0.361 \text{WBC} - 0.361}}{e^{1.737 - 0.361 \text{WBC}}} = e^{-0.361} \end{aligned}$$

a one-unit increase in WBC is associated with a change in odds of dengue by a factor of $e^{-0.361} \approx 0.7$