

Logistic regression

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)

Research goal: Predict dengue status using diagnostic measurements

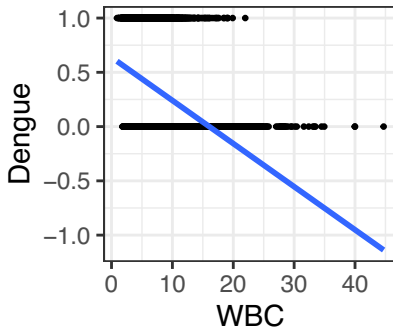
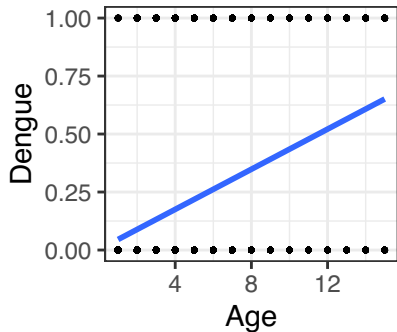
Last time: initial attempt

What if we try a linear regression model?

Y_i = dengue status of i th patient

$$Y_i = \beta_0 + \beta_1 WBC_i + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

Don't fit linear regression with a binary response



Last time: rewriting the linear regression model

$$Y_i | WBC_i \sim N(\mu_i, \sigma^2) \quad (\text{random component})$$

$$\mu_i = \beta_0 + \beta_1 WBC_i \quad (\text{systematic component})$$

describe of distribution
of $Y_i | WBC_i$

↑
specifies new parameters
of the distribution
depend on WBC

Second attempt

$$Y_i | WBC_i \sim \text{Bernoulli}(p_i) \quad p_i = \mathbb{P}(Y_i = 1 | WBC_i)$$

$$p_i = \beta_0 + \beta_1 WBC_i$$

Are there still any potential issues with this approach?

Problem : $p_i \in [0, 1]$

but $\beta_0 + \beta_1 WBC_i \in (-\infty, \infty)$
(potentially)

Fixing the issue: logistic regression

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

(random component)

link function $\nearrow g(p_i) = \beta_0 + \beta_1 WBC_i$

(systematic component)

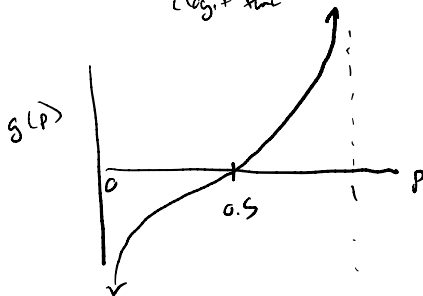
where $g : (0, 1) \rightarrow \mathbb{R}$ is unbounded.

Usual choice: $g(p_i) = \log\left(\frac{p_i}{1-p_i}\right)$
(logit function)

$$p_i \in (0, 1)$$

$$\frac{p_i}{1-p_i} = \frac{\text{odds}}{\in (0, \infty)}$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \log \text{odds} \in (-\infty, \infty)$$



Logistic regression model

$$Y_i | WBC_i \sim \text{Bernoulli}(p_i) \quad \leftarrow \text{capturing randomness about } Y_i$$

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

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.

Linear model :

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

\Rightarrow

$$Y_i | X_i \sim N(\mu_i, \sigma^2)$$

$$\mu_i = \beta_0 + \beta_1 X_i$$

Fitting the logistic regression model

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

"generalized
linear model"



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

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

family of specifies distribution
of the response

Fitting the logistic regression model

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

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

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

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

```
##
```

```
## Call:
```

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

```
##
```

```
## 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

Activity

Work on the activity (handout) in groups, then we will discuss as a class.

Activity

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

What is the predicted *odds* of dengue for a patient with a white blood cell count of 15?

$$\begin{aligned} \log \hat{\text{odds}} &= 1.737 - 0.361(15) \\ &= -3.678 \end{aligned}$$

$$\hat{\text{odds}} = e^{-3.678} = 0.0253$$

Activity

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

What is the predicted *probability* of dengue for a patient with a WBC of 15?

$$\text{odds} = \frac{p}{1-p} \quad \Leftrightarrow \quad p = \frac{\text{odds}}{1 + \text{odds}}$$

$$= \frac{e^{\log \text{odds}}}{1 + e^{\log \text{odds}}}$$

$$\hat{p}_i = \frac{e^{-3.678}}{1 + e^{-3.678}} \approx 0.025$$

Interpretation: Activity 2

Work on the activity (handout) in groups, then we will discuss as a class.

Interpretation

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

Are patients with a higher WBC more or less likely to have dengue?

less likely (negative slope)
 $\hat{\beta}_1 < 0$

Interpretation

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

What is the change in *log odds* associated with a unit increase in WBC?

- 0.361

Interpretation

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

What is the change in *odds* associated with a unit increase in WBC?

$$\begin{aligned} \frac{\text{odds when WBC} = x+1}{\text{odds when WBC} = x} &= \frac{e^{1.737 - 0.361(x+1)}}{e^{1.737 - 0.361x}} \\ &= e^{-0.361} \\ &= 0.697 \end{aligned}$$

A one unit increase in WBC is associated with a change in the odds of dengue by a factor of 0.697 (30.3% decrease)

Coefficient interpretation

$$\log \left(\frac{\hat{p}_i}{1 - \hat{p}_i} \right) = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

Fitting a *logistic* regression model?

Linear regression: minimize $\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_k X_{ik})^2$

Question: Should we minimize a similar sum of squares for a *logistic* regression model?