

Fitting logistic regression models

where did the ε_i go?

$$\gamma Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \text{random}$$

$$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma_\varepsilon^2)$$

$$Y_i | X_i \sim N(\beta_0 + \beta_1 X_i, \sigma_\varepsilon^2)$$

γY_i is not iid b/c μ_i (or p_i) changes, so not identically distributed



$$\left. \begin{array}{l} Y_i | X_i \sim N(\mu_i, \sigma_\varepsilon^2) \\ \mu_i = \beta_0 + \beta_1 X_i \end{array} \right\} \begin{array}{l} \text{(random component)} \\ \text{(Systematic component)} \end{array}$$

not random

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

(random component)

$$\log\left(\frac{p_i}{1-p_i}\right) = \underbrace{\beta_0 + \beta_1 X_i}_{\in (-\infty, \infty)} \quad \text{(Systematic component)}$$

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

$$Y_i = p_i + \varepsilon_i$$

$$\varepsilon_i = \begin{cases} 1 - p_i & \text{w/prob } p_i \\ -p_i & \text{w/prob } 1 - p_i \end{cases}$$

Motivating example: 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)

Last time: Logistic regression model

Y_i = dengue status (0 = negative, 1 = positive)

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

random

link function \rightarrow

$$\log\left(\frac{p_i}{1 - p_i}\right) = \underbrace{\beta_0 + \beta_1 WBC_i}_{\text{linear}} \quad \text{systematic}$$

We get n observations $(WBC_1, Y_1), \dots, (WBC_n, Y_n)$. Want estimates $\hat{\beta}_0, \hat{\beta}_1$

1 unit increase in WBC is associated w/ a change in odds of dengue by a factor of 0.697

Last time: Logistic regression model

$Y_i = \text{dengue status } (0 = \text{no}, 1 = \text{yes}) \quad Y_i \sim \text{Bernoulli}(p_i)$

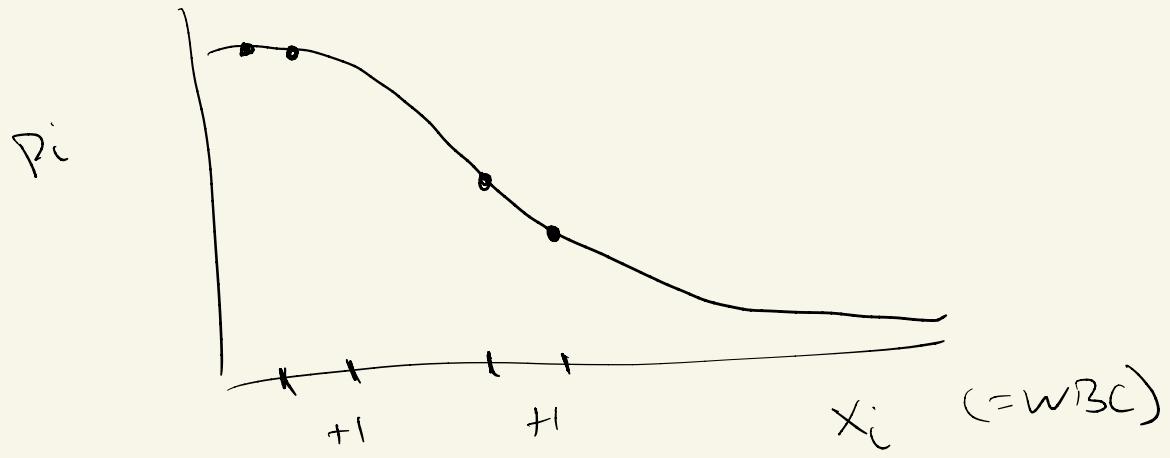
$$\log\left(\frac{\hat{p}_i}{1 - \hat{p}_i}\right) = 1.737 - 0.361 WBC_i$$

How should we interpret the slope -0.361?

log odds: 1 unit increase in WBC is associated w/ a decrease of 0.361 in log odds of dengue

$$\text{odds} : \text{odds} = \exp^{\{1.737 - 0.361 WBC_i\}}$$

$$\begin{aligned} WBC &\rightarrow WBC + 1 \\ \exp^{\{1.737 - 0.361 WBC\}} &\rightarrow \exp^{\{1.737 - 0.361 WBC - 0.361\}} \\ \text{odds ratio} &= \frac{e}{e^{\{1.737 - 0.361 WBC\}}} \\ &= e^{-0.361} \\ &\approx 0.697 \end{aligned}$$



$$\hat{p}_i \in [0, 1]$$

Getting probabilities

$Y_i = \text{dengue status } (0 = \text{no}, 1 = \text{yes}) \quad Y_i \sim \text{Bernoulli}(p_i)$

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

How do I calculate estimated probabilities \hat{p}_i ?

$$\frac{\hat{p}_i}{1 - \hat{p}_i} = e^{1.737 - 0.361 WBC_i}$$

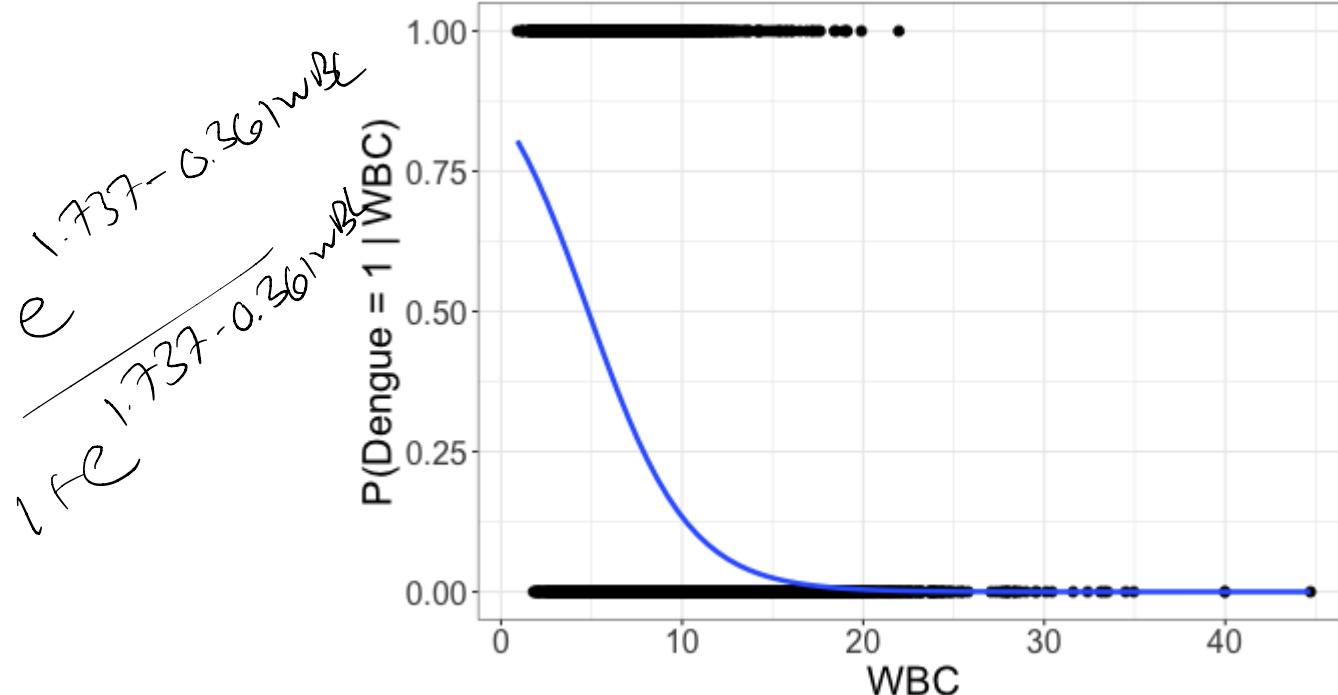
$$\hat{p}_i = e^{1.737 - 0.361 WBC_i} (1 - \hat{p}_i)^{1.737 - 0.361 WBC_i}$$

$$\hat{p}_i(1 + e^{1.737 - 0.361 WBC_i}) = e^{1.737 - 0.361 WBC_i}$$

$$\Rightarrow \hat{p}_i = \frac{e^{1.737 - 0.361 WBC_i}}{1 + e^{1.737 - 0.361 WBC_i}} = \frac{\text{odds}}{1 + \text{odds}}$$

Curve

Plotting the fitted model for dengue data

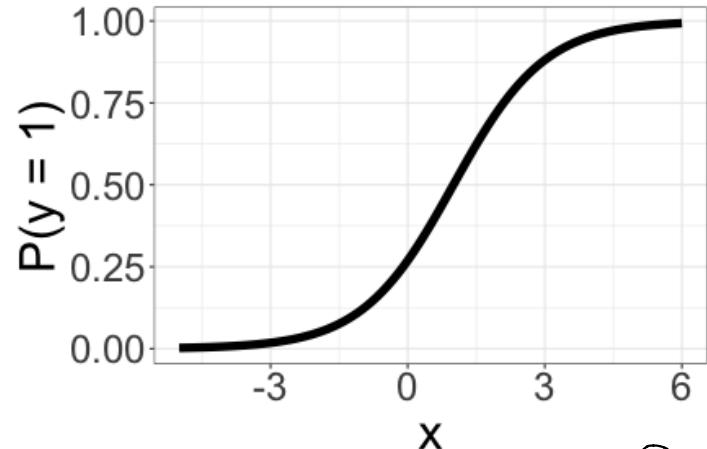
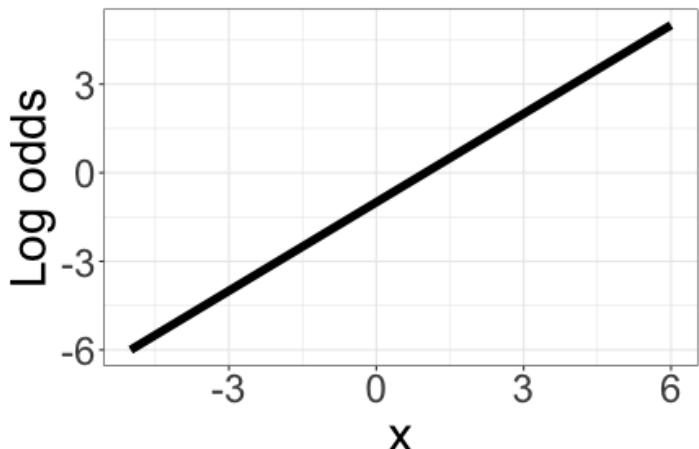


Shape of the regression curve

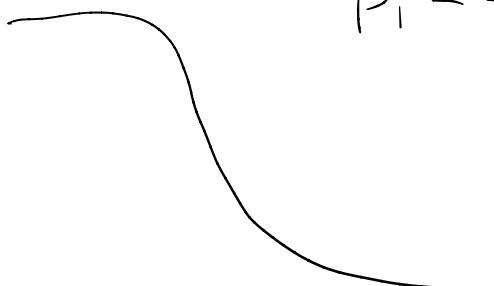
$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 X_i$$

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

$\beta_1 > 0$



$\beta_1 < 0$



Shape of the regression curve

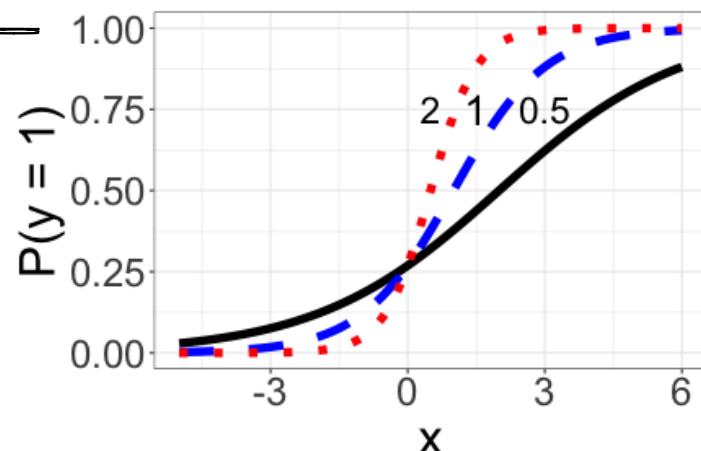
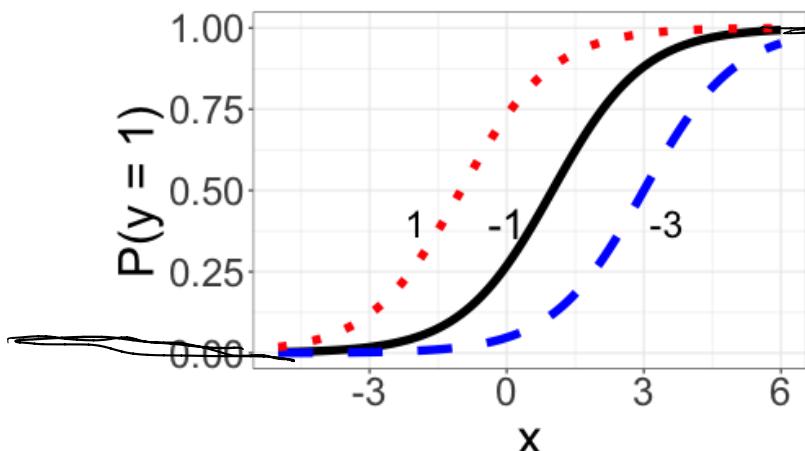
How does the shape of the fitted logistic regression depend on β_0 and β_1 ?

$$p_i = \frac{\exp\{\beta_0 + X_i\}}{1 + \exp\{\beta_0 + X_i\}}$$

for $\beta_0 = -3, -1, 1$

$$p_i = \frac{\exp\{-1 + \beta_1 X_i\}}{1 + \exp\{-1 + \beta_1 X_i\}}$$

for $\beta_1 = 0.5, 1, 2$



Fitting logistic regression in R

vs. lm

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

```
summary(m1)
```

Distribution of response

...

	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	***
## ---					
## Signif. codes:	0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1				
##	(ignores for now)				
## (Dispersion parameter for binomial family taken to be 1)					
##					
## Null deviance:	6955.8	on 5719	degrees of freedom		
## Residual deviance:	5529.8	on 5718	degrees of freedom		
## AIC:	5533.8	instead of R^2			
##					
## Number of Fisher Scoring iterations:	5				

not t!

Recap: ways of fitting a *linear* regression model

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \cdots + \beta_k X_{i,k} + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

How do we fit this linear regression model? That is, how do we estimate

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

Discuss with your neighbor for 2--3 minutes.

- 1) minimize SSE
- 2) Projection ($\Rightarrow (X^T X)^{-1} X^T Y$)
- 3) MLE

Method 1: Minimize SSE

$$SSE = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2$$

squared residuals

$$\frac{\partial SSE}{\partial \beta_0} \quad \begin{matrix} \text{set} \\ = 0 \end{matrix}$$

$$\frac{\partial SSE}{\partial \beta_1} \quad \begin{matrix} \text{set} \\ = 0 \end{matrix}$$

$$\vdots \quad \vdots \quad \begin{matrix} \text{set} \\ = 0 \end{matrix}$$

$$\frac{\partial SSE}{\partial \beta_k} \quad \begin{matrix} \text{set} \\ = 0 \end{matrix}$$

$n+1$ equations
 $n+1$ unknowns
⇒ solve

Method 2: Projection argument

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \in \mathbb{R}^n$$

$$X = \begin{pmatrix} 1 & X_{11} & \dots \\ 1 & X_{21} & \dots \\ \vdots & \vdots & \ddots \\ 1 & X_{n1} & \dots \end{pmatrix}$$

$$\beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}$$

$$\hat{Y} = X \hat{\beta}$$

want \hat{Y} "close" to Y

$$\Rightarrow \|Y - \hat{Y}\| \text{ to be small}$$
$$\|Y - \hat{Y}\| = \left(\sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_k X_{ik})^2 \right)^{1/2}$$

\Leftrightarrow minimizing SSE!

Method 3: Maximizing likelihood

$$Y_i \sim N(\beta_0 + \beta_1 X_{i1} + \dots + \beta_n X_{in}, \sigma^2)$$

$$L(\beta_0, \dots, \beta_n, \sigma^2) = \prod_{i=1}^n f(Y_i; \beta_0, \beta_1, \dots, \beta_n, \sigma^2)$$

$$= (2\pi)^{-\frac{n}{2}} \sigma^{-n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_n X_{in})^2 \right\}$$

SSE!

maximizing likelihood \Leftrightarrow minimizing SSE
(for normal data)

Summary: three ways of fitting linear regression models

- + Minimize SSE, via derivatives of

$$\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i,1} - \cdots - \beta_k X_{i,k})^2$$

- + Minimize $\|\hat{Y}\|$ (equivalent to minimizing SSE)
- + Maximize likelihood (for *normal* data, equivalent to minimizing SSE)

Which of these three methods, if any, is appropriate for fitting a logistic regression model? Do any changes need to be made for the logistic regression setting?

Discuss with your neighbor for 2--3 minutes.

Maximum likelihood for logistic regression

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

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 X_{i,1} + \cdots + \beta_k X_{i,k}$$

Suppose we observe independent samples $(X_1, Y_1), \dots, (X_n, Y_n)$. Write down the likelihood function

$$L(\beta) = \prod_{i=1}^n f(Y_i; \beta)$$

for the logistic regression problem. Take 2--3 minutes, then we will discuss as a class.

Maximum likelihood for logistic regression

$$L(\beta) =$$

I want to choose β to maximize $L(\beta)$. What are the usual steps to take?

Initial attempt at maximizing likelihood

$$L(\beta) = \prod_{i=1}^n p_i^{Y_i} (1 - p_i)^{1-Y_i}$$

$$\ell(\beta) =$$

Iterative methods for maximizing likelihood

Fisher scoring

Fisher scoring for logistic regression