

Modeling a Binary Response

Logistic regression – Stat 230

Framingham Heart Study

- Long-term study of cardiovascular disease (heart attack, stroke, and other heart problems)
- First study of cardiovascular disease to follow initially healthy people over a long period of time
- Identified the idea that body measurements (blood pressure, cholesterol, etc) could predict cardiovascular disease.
- Participants from a population of free living (not hospitalized) subjects in the community of Framingham, Mass.
- Started in 1948, participants have been examined every other year since the inception
- In 3rd generation of participants!

Modeling overview

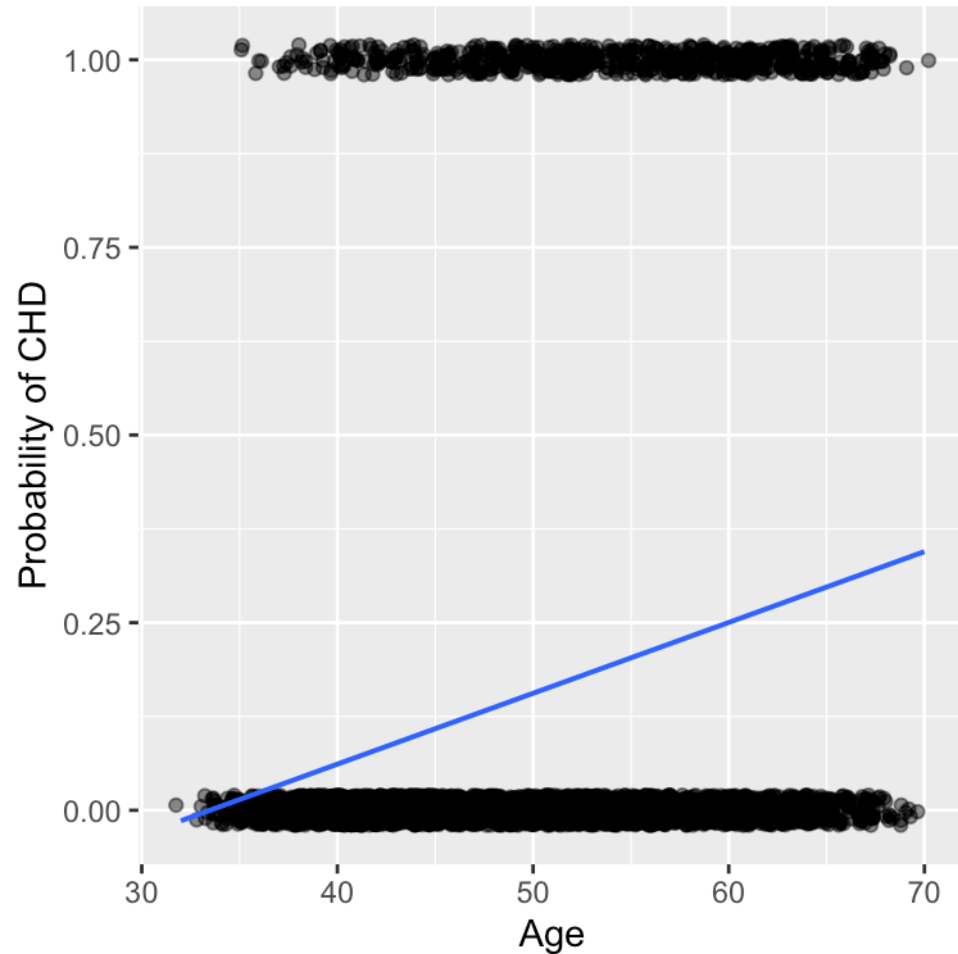
Goal: Determine whether participant experienced coronary heart disease (CHD) in 10-year window after their exam

- $Y = \text{CHD}$ (0 = no, 1 = yes)
- X = participant's age, sex, total cholesterol, and systolic blood pressure

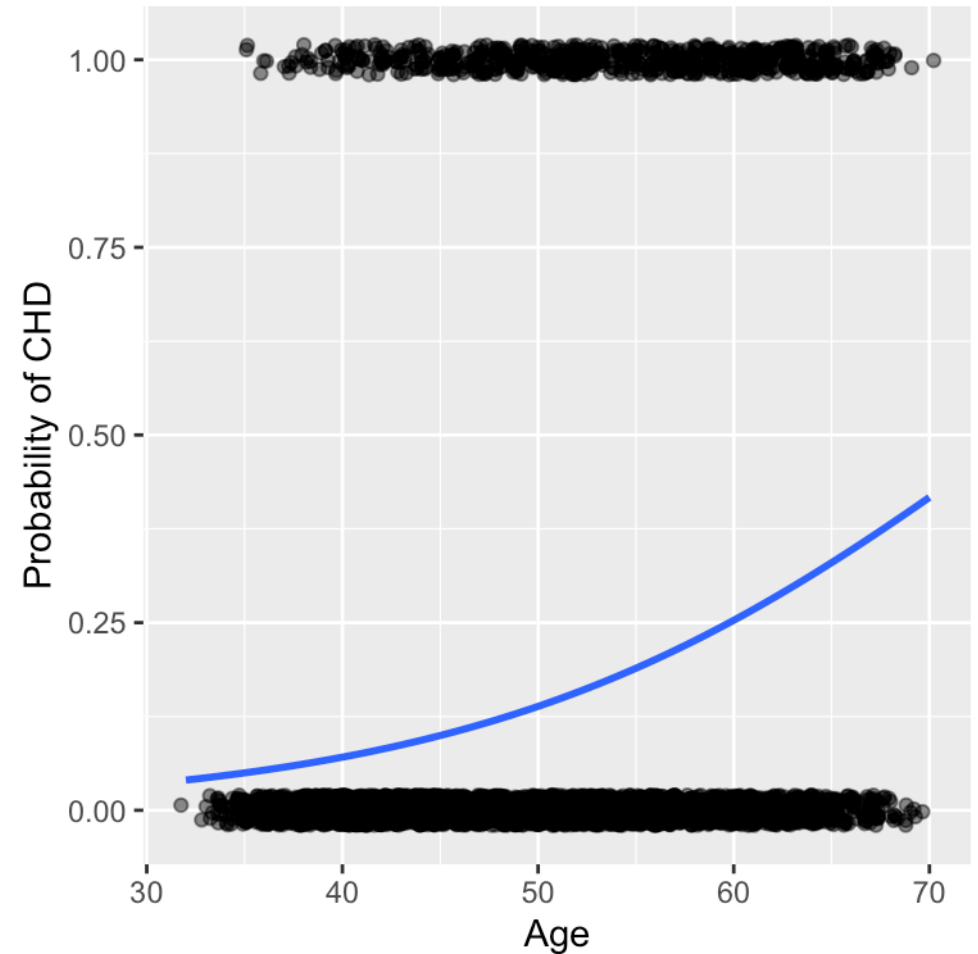
Strategy: model the probability of CHD given these factors

How can we model a binary response?

Linear regression



Logistic regression



Generalized linear model

Random component:

What distribution does the response variable follow?

Linear predictor (systematic component):

What function of the predictors should we use?

Link function:

How can we relate the expectation of Y and our linear predictor?

Random component

What distribution does the response variable follow?

Possible values: $Y = 0$ or 1

Bernoulli distribution

$$P(Y = y) = \pi^y(1 - \pi)^{1-y}, \quad y = 0, 1$$

Linear predictor

What function of the predictors should we use?

$$\eta = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

Link function

How can we relate the expectation of Y and our linear predictor?

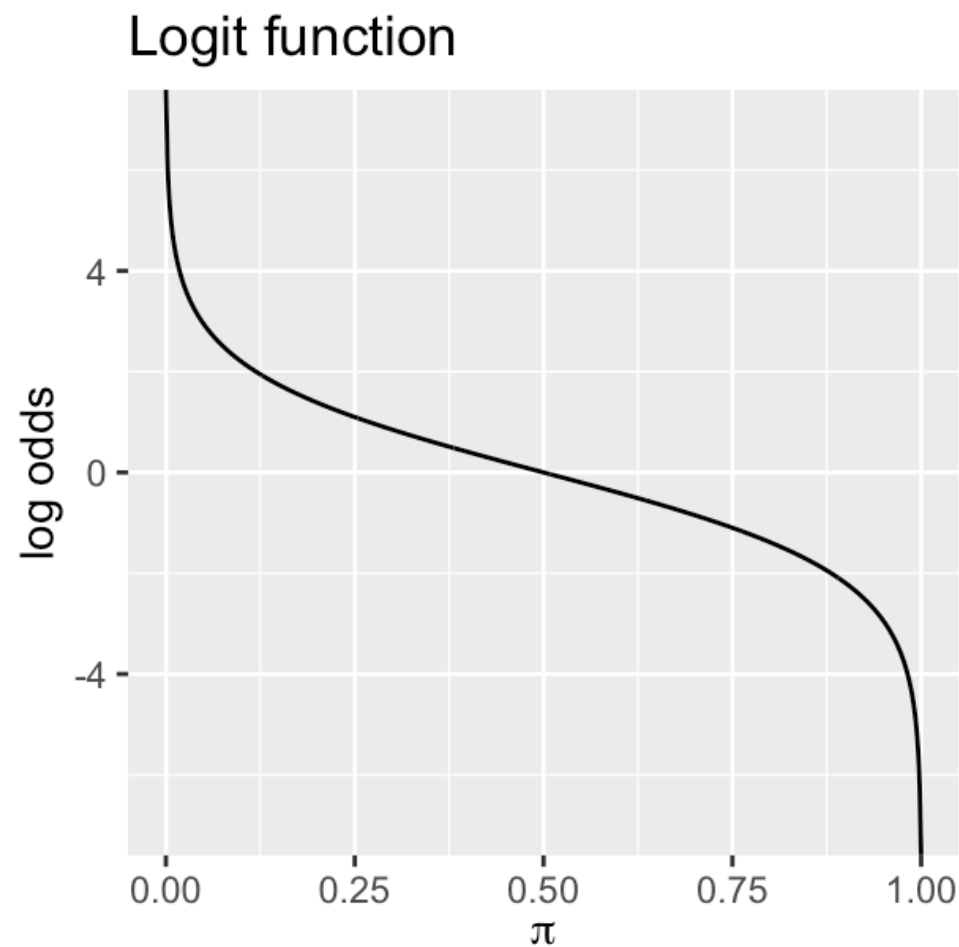
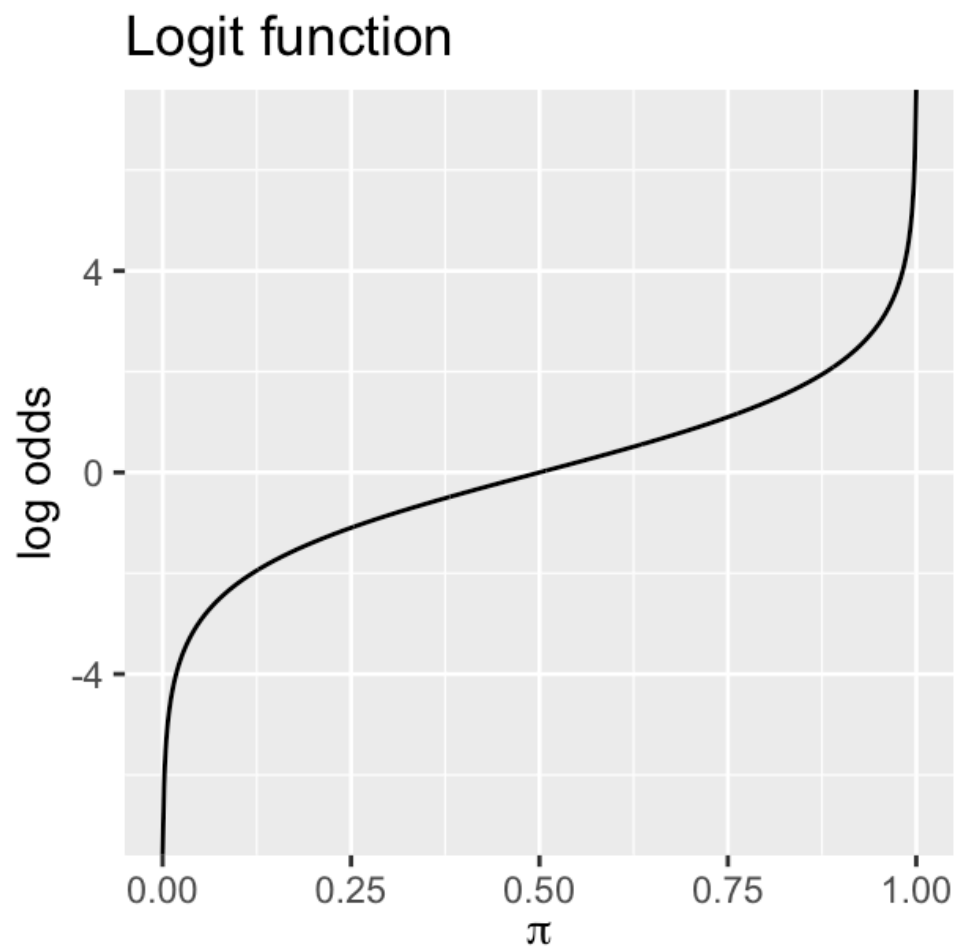
$$E(Y) = \pi \longrightarrow \text{between } 0 \text{ and } 1$$

$$\eta = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p \longrightarrow \text{between } -\infty \text{ and } \infty$$

What function takes numbers between 0 and 1 as input and maps them to the real line?

The logit function

$$\text{logit}(\pi) = \log\left(\frac{\pi}{1 - \pi}\right)$$



Link function

How can we relate the expectation of Y and our linear predictor?

$$\text{logit}(\pi) = \log\left(\frac{\pi}{1 - \pi}\right) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p = \eta$$

Example: A quantitative predictor

Let's use total cholesterol as a predictor of CHD

$$\hat{\eta}_i = \log \left(\frac{\hat{\pi}(\text{totChol}_i)}{1 - \hat{\pi}(\text{totChol}_i)} \right) = -2.894 + 0.005 \text{totChol}_i$$

$\hat{\beta}_1 = 0.005$ is the expected change in the log odds associated with a 1-unit increase in total cholesterol...

We don't think in terms of log odds, so we need more interpretable quantities!

Log odds \rightarrow odds

Log odds:

$$\hat{\eta} = \log \left(\frac{\hat{\pi}(X)}{1 - \hat{\pi}(X)} \right) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p$$

Odds of success:

$$\hat{\omega} = \frac{\hat{\pi}(X)}{1 - \hat{\pi}(X)} = e^{\hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p}$$

Example: A quantitative predictor

$$\hat{\eta}_i = \log\left(\frac{\hat{\pi}(\text{totChol})}{1 - \hat{\pi}(\text{totChol})}\right) = -2.894 + 0.005\text{totChol}_i$$

How do we interpret $e^{\hat{\beta}_1} = e^{0.005} \approx 1.005$??

$$\frac{\hat{\pi}(\text{totChol})}{1 - \hat{\pi}(\text{totChol})} = e^{-2.894 + 0.005\text{totChol}_i} = e^{-2.894} \cdot e^{0.005\text{totChol}_i}$$

So if **totChol** increases by one unit...

$$\begin{aligned}\frac{\hat{\pi}(\text{totChol} + 1)/\{1 - \hat{\pi}(\text{totChol} + 1)\}}{\hat{\pi}(\text{totChol})/\{1 - \hat{\pi}(\text{totChol})\}} &= \frac{e^{-2.894} \cdot e^{0.005\text{totChol}+1}}{e^{-2.894} \cdot e^{0.005\text{totChol}}} \\ &= \frac{e^{-2.894} \cdot e^{0.005\text{totChol}} \cdot e^{0.005(1)}}{e^{-2.894} \cdot e^{0.005\text{totChol}}} \\ &= e^{0.005(1)} \approx 1.005\end{aligned}$$

we expect the odds of CHD to increase by a factor of 1.005

Estimating probabilities

What's the probability that a person with total cholesterol of 240 mg/dL will develop CHD in the next 10 years?

$$\hat{\eta}(240) = \log \left(\frac{\hat{\pi}(240)}{1 - \hat{\pi}(240)} \right) = -2.894 + 0.005(240) = -1.694$$

How do we get from the estimated log odds to a probability?

Log odds \rightarrow probability

Log odds:

$$\hat{\eta} = \log \left(\frac{\hat{\pi}(X)}{1 - \hat{\pi}(X)} \right) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p$$

Odds of success:

$$\hat{\omega} = \frac{\hat{\pi}(X)}{1 - \hat{\pi}(X)} = e^{\hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p}$$

Probability of success:

$$\hat{\pi}(X) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 X_1 + \cdots + \hat{\beta}_p X_p}}$$

Estimating probabilities

What's the probability that a person with total cholesterol of 240 mg/dL will develop CHD in the next 10 years?

$$\hat{\eta}(240) = \log \left(\frac{\hat{\pi}(240)}{1 - \hat{\pi}(240)} \right) = -2.894 + 0.005(240) = -1.694$$

How do we get from the estimated log odds to a probability?

$$\hat{\pi}(240) = \frac{e^{-1.694}}{1 + e^{-1.694}} \approx 0.155$$

Fitted logistic regression model

$\text{logistic}(\hat{\eta}) = \frac{\exp(\hat{\eta})}{1 + \exp(\hat{\eta})}$ is the fitted curve on the probability scale

