STAT 608, Spring 2022 - Assignment 6
SOLUTIONS

1. For logistic regression with one predictor, we use the model

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

(a) Show that solving for the probability of success for a given value of the predictor, $\theta(x)$, gives

$$\theta(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

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

$$\frac{\theta(x)}{1 - \theta(x)} = \exp(\beta_0 + \beta_1 x)$$

$$\theta(x) = \exp(\beta_0 + \beta_1 x) - \theta(x)\exp(\beta_0 + \beta_1 x)$$

$$\theta(x) + \theta(x)\exp(\beta_0 + \beta_1 x) = \exp(\beta_0 + \beta_1 x)$$

$$\theta(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

(b) and

$$\theta(x) = \frac{1}{1 + \exp(-\{\beta_0 + \beta_1 x\})}$$

$$\theta(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$
$$= \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}\left[\frac{\exp(-(\beta_0 + \beta_1 x))}{\exp(-(\beta_0 + \beta_1 x))}\right]$$
$$= \frac{1}{\exp(-(\beta_0 + \beta_1 x)) + 1}$$
$$= \frac{1}{1 + \exp(-(\beta_0 + \beta_1 x))}$$

2. On page 285 of the text, it says "When $X$ is a dummy variable, it can be shown that the log odds are also a linear function of $x$." Suppose that $X$ is a dummy variable, taking the value 1 with probability $\pi_j$, $j = 0, 1$, conditional on $Y = 0, 1$.

(a) Show that the log odds are a linear function of $x$.

FIRST, DEFINE THE BERNOULLI PROBABILITY $P(X|Y = j) = \pi_j^x(1 - \pi_j)^{1-x}$. THEN NOTE THAT

$$\frac{\theta(x)}{1 - \theta(x)} = \frac{P(Y = 1|X)}{P(Y = 0|X)} = \frac{P(Y = 1, X = x)P(X = x)}{P(Y = 0, X = x)P(X = x)} = \frac{P(Y = 1)P(X = x|Y = 1)}{P(Y = 0)P(X = x|Y = 0)}$$

FINALLY, WE HAVE

$$\log\left(\frac{\theta(x)}{1 - \theta(x)}\right) = \log\left(\frac{P(Y = 1)}{P(Y = 0)}\right) + \log\left(\frac{P(X = x|Y = 1)}{P(X = x|Y = 0)}\right)$$

$$= \log\left(\frac{P(Y = 1)}{P(Y = 0)}\right) + \log\left(\frac{\pi_1^x(1 - \pi_1)^{1-x}}{\pi_0^x(1 - \pi_0)^{1-x}}\right)$$

$$= \log\left(\frac{P(Y = 1)}{P(Y = 0)}\right) + x\log\left(\frac{\pi_1}{\pi_0}\right) + (1 - x)\log\left(\frac{1 - \pi_1}{1 - \pi_0}\right)$$

$$= \log\left(\frac{P(Y = 1)}{P(Y = 0)}\right) + \log\left(\frac{1 - \pi_1}{1 - \pi_0}\right) + x\left(\log\left(\frac{\pi_1}{\pi_0}\right) - \log\left(\frac{1 - \pi_1}{1 - \pi_0}\right)\right)$$

$$= \log\left(\frac{P(Y = 1)P(X = 0|Y = 1)}{P(Y = 0)P(X = 0|Y = 0)}\right) + x\log\left(\frac{\pi_1/(1 - \pi_1)}{pi_0/(1 - \pi_0)}\right)$$

$$= a + bx$$

(b) Define the slope and intercept for the linear function.

THE INTERCEPT IS

$$\log\left(\frac{P(Y = 1)}{P(Y = 0)}\right) + \log\left(\frac{1 - \pi_1}{1 - \pi_0}\right)$$

AND THE SLOPE IS

$$\log\left(\frac{\pi_1/(1 - \pi_1)}{\pi_0/(1 - \pi_0)}\right)$$

3. On page 284 of the text, the author quotes Cook and Weisberg: "When conducting a binary regression with a skewed predictor, it is often easiest to assess the need for $x$ and $\log(x)$ by including them both in the model so that their relative contributions can be assessed directly." Show that indeed the log odds are a function of $x$ and $\log(x)$ for the gamma distribution.

THE TEXT GIVES THE FOLLOWING, SO WE HANDLE ONLY THE SECOND TERM IN THE SUM, AS THE FIRST TERM IS A CONSTANT WITH RESPECT TO $x$:

$$\log\left(\frac{\theta(x)}{1 - \theta(x)}\right) = \log\left(\frac{P(Y = 1)}{P(Y = 0)}\right) + \log\left(\frac{f(x|Y = 1)}{f(x|Y = 0)}\right)$$

WE'LL USE THE PARAMETRIZATION OF THE GAMMA DISTRIBUTION AS FOLLOWS:

$$f(x|\alpha, \beta) = \frac{x^{\alpha-1}e^{-x/\beta}}{\Gamma(\alpha)\beta^\alpha}$$

THEN WE CAN START REARRANGING:

$$\log\left(\frac{f(x|Y=1)}{f(x|Y=0)}\right) = \log\left(\frac{x^{\alpha_1-1}e^{-x/\beta_1}\Gamma(\alpha_0)\beta_0^{\alpha_0}}{x^{\alpha_0-1}e^{-x/\beta_0}\Gamma(\alpha_1)\beta_1^{\alpha_1}}\right)$$

$$= \log\left(\frac{\Gamma(\alpha_0)\beta_0^{\alpha_0}}{\Gamma(\alpha_1)\beta_1^{\alpha_1}}\right) + (\alpha_1-\alpha_0)\log(x) + \left(\frac{1}{\beta_0}+\frac{1}{\beta_1}\right)x$$

FINALLY, WE SEE THAT THE LOG ODDS IS A FUNCTION OF $x$ AND $\log(x)$ WHEN $x$ HAS THE GAMMA DISTRIBUTION.

4. Chapter 8, Question 4

   (a) NO, MODEL (8.6) IS NOT A VALID MODEL: WE SEE FROM THE MARGINAL MODEL PLOTS THAT OUR MODEL DOES NOT MATCH A NONPARAMETRIC FIT TO THE DATA.

   (b) THE REAL ANSWER HERE IS THAT THE PREDICTORS MAY VERY WELL HAVE NONLINEAR ASSOCIATIONS AND WE MAY BE MISSING SOME IMPORTANT PREDICTORS LIKE DIET AND EXERCISE; HOWEVER, BASED ON THE KERNEL DENSITY ESTIMATES, WE FIRST OBSERVE THAT $x_1$ AND $x_4$ ARE BOTH RIGHT-SKEWED. AS SHOWN ABOVE, THIS MAY INDICATE THAT ADDING THE LOG TRANSFORMATIONS OF BOTH OF THESE VARIABLES MAY BE IMPORTANT. TO A SMALLER DEGREE, IT APPEARS VARIABLE $x_1$ HAS UNEQUAL VARIANCES WHEN HEART DISEASE = YES AND NO, WHICH MAY INDICATE ADDING A QUADRATIC TERM IF WE THINK $x_1$ IS APPROXIMATELY NORMALLY DISTRIBUTED. THE FACT THAT THE PLOT IS MUCH WIDER THAN IT IS NARROW MAY BE EXAGGERATING THE SKEWNESS TO OUR EYES.

   (c) THIS IS A MUCH-IMPROVED MODEL. WE SEEM TO STILL BE HAVING TROUBLE MODELING THE RELATIONSHIP BETWEEN $x_1$ AND THE ODDS OF A HEART ATTACK AT VERY SMALL VALUES OF $x_1$, BUT OTHERWISE, THE MODEL SEEMS TO MATCH A NONPARAMETRIC FIT PRETTY WELL. WE DON'T HAVE THE ABILITY TO CONDUCT A GOODNESS OF FIT TEST FOR THESE BINARY DATA.

   (d) HOLDING SYSTOLIC BLOOD PRESSURE, CHOLESTEROL, OBESITY, AND AGE CONSTANT, WHEN A PATIENT HAS A FAMILY HISTORY OF HEART DISEASE, OUR MODEL PREDICTS THEIR ODDS OF A HEART ATTACK TO BE $e^{0.941056} = 2.56$ TIMES AS LARGE AS WHEN A PATIENT DOES NOT HAVE A FAMILY HISTORY OF HEART DISEASE.