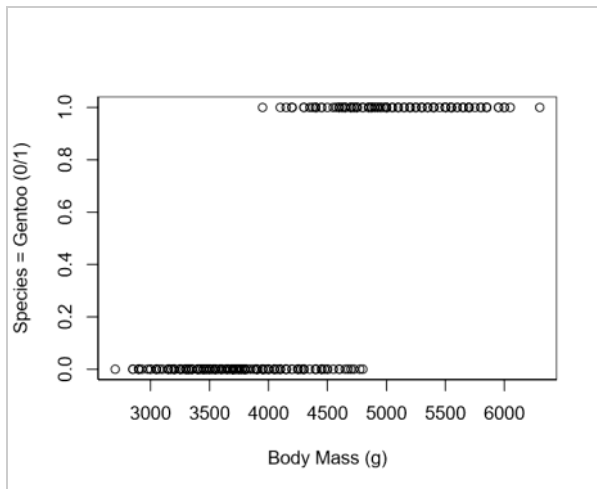


Logistic Regression

Previously

Regression: Describe average value of a continuous response variable Y from one or more variables X_1, X_2, X_3, \dots numeric (discrete or continuous)

What if the variable Y were categorical rather than numeric? (Note: we focus on binary variables)



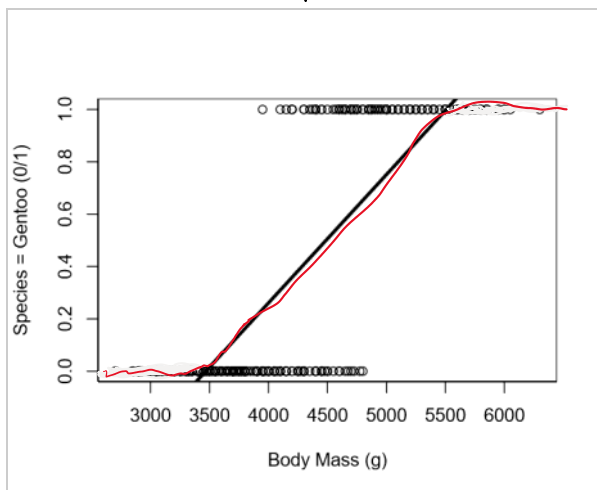
Body mass as a predictor of Species being Gentoo.

$$Y = \{0, 0, 1, 1, 0, \dots\}$$

Chinstrap or Adelie
Gentoo

What if we just change nothing?

$$\hat{y} = b_0 + b_1 \cdot BM \quad \leftarrow \text{Let's fit this model.}$$



This is our fit.

How can we fix it?

$\hat{y} > 1 \rightarrow$ replace with 1

$\hat{y} < 0 \rightarrow$ replace with 0

Not a statistically motivated solution \nearrow

This does not work.

The problem is $\hat{y} = b_0 + b_1 \cdot BM$ Takes values on real number line
 $\mathbb{R} = (-\infty, +\infty)$

The problem is $y = b_0 + b_1 BM$ $R = (-\infty, +\infty)$

mean for a given BM \nwarrow Takes values of 0 or 1.

We want to replace the mean value \hat{y} with a prob (Specifically $P(\text{Species} = \text{Gentoo})$)

$$p = b_0 + b_1 BM$$

\nwarrow $[0, 1]$ \nwarrow $(-\infty, \infty)$

2 options

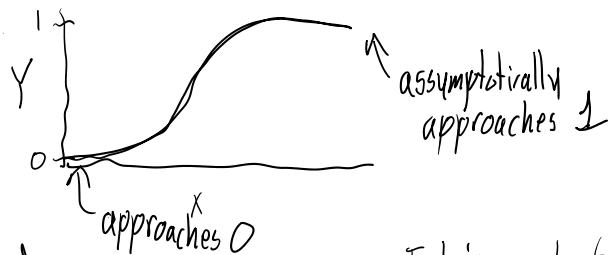
Both sides of equation are in range $(-\infty, \infty)$

or
Both sides in range $(0, 1)$
(We will do both)

What we will fit is

$$\text{Function}(p) = b_0 + b_1 BM$$

Sigmoid



We need this
function to map values of
 p from $(0, 1) \rightarrow (-\infty, \infty)$
 \nwarrow $p = P(\text{Gentoo})$

$$\log\left(\frac{p}{1-p}\right) = b_0 + b_1 BM$$

Sigmoid

$$y = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

Euler's constant
 ≈ 2.71

Note

$$\log(a) = b \quad e^{a+b} = e^a \cdot e^b$$

$$a = e^b \quad e^{a-b} = e^a / e^b$$

What is $\frac{p}{1-p}$?

This is the odds.

p is a probability which
measures the number of
times an event will occur
given all possible events

$$P(\text{Coin} = \text{heads}) = 1 \text{ in } 2 = 1/2$$

$$P(\text{Die Roll} = 3) = 1 \text{ in } 6 = 1/6$$

If we solve for x ...

$$\log\left(\frac{y}{1-y}\right) = x$$

\nwarrow natural loge

$p/1-p$ is an odds, which
measures the number of
times an event occurs
compared to the number
of times it doesn't occur.

$$\text{Odds}(\text{Coin} = \text{heads}) = 1 \text{ to } 1 = 1$$

$$\text{Odds}(\text{Die Roll} = 3) = 1 \text{ to } 5 = 1/5$$

$$P(\text{Die Roll} = 3) = 1 \text{ in } 6 = 1/6$$

↖ [0,1]

$$\text{Odds}(\text{Coin-flip}) = 1 \text{ to } 1 = 1$$

$$\text{Odds}(\text{Pie Roll} = 3) = 1 \text{ to } 5 = 1/5$$

↖ $(-\infty, \infty)$

We define the log-odds using a linear equation

$$\log\left(\frac{p}{1-p}\right) = \ln\left(\frac{p}{1-p}\right) = b_0 + b_1 \text{ BM}$$

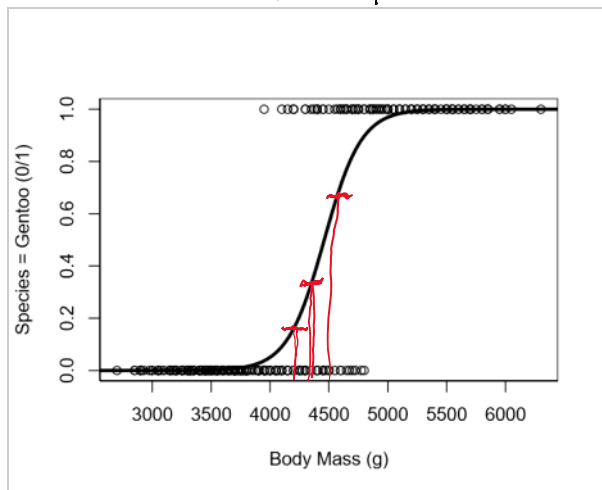
units of grams

Using R: $\hat{b}_0 = -28.41806$
 $\hat{b}_1 = 0.00637 \frac{1}{g}$

$$p = \frac{1}{1 + e^{-(b_0 + b_1 x)}}$$

$$= \frac{e^{b_0 + b_1 x}}{1 + e^{b_0 + b_1 x}}$$

$$p = \frac{1}{1 + \exp(-(-28.418 + 0.00637 \cdot \text{BM}))}$$



Interpretation of terms is complex, but easiest if we look at odds.

$$\log\left(\frac{p}{1-p}\right) = b_0 + b_1 \text{ BM}$$

$$\frac{p}{1-p} = e^{b_0 + b_1 \text{ BM}}$$

$$= e^{b_0} \times e^{b_1 \cdot \text{BM}}$$

Before 1 unit increase in BM would result in a \hat{b}_1 increase in Y on average.

Here 1 unit increase in BM increases the odds multiplicatively by $e^{\hat{b}_1}$

A one unit increase in body mass is associated with a 0.6% increase in odds.

Penguin BM = 4000
and
BM = 4001

$$\frac{e^a}{e^b} = e^{a-b}$$

$$\frac{\left(\frac{p}{1-p}\right)_{4001}}{\left(\frac{p}{1-p}\right)_{4000}} = e^{\hat{b}_0 + \hat{b}_1 \cdot 4001 - (\hat{b}_0 + \hat{b}_1 \cdot 4000)}$$

$$= e^{\hat{b}_1}$$

$$e^{0.00637} = 1.00639$$

BM is correlated

in odds.

C

- 1.00629

A ten unit increase in BM is associated
with a 1.0658 times increase in odds
(6.58%)

A 100 unit increase in BM ----
..... 89.1% increase (1.89 times)

1000 unit increase ----

↑
1 kg 58300% increase