

Lecture 4: Logistic Regression

*Instructor: Ryan Brill**Scribe: Jonathan Pipping*

4.1 Motivating Example: Sinking Putts

4.1.1 Problem Setup

Question: Predict the **probability** that a putt is sunk as a function of the distance to the hole.

Data: 5,988 putts from Columbia including distance to the hole and whether the putt was sunk or not. Each row in the dataset is a single putt and includes the following variables:

- i : index of the i^{th} putt in our dataset
- y_i : 1 if the i^{th} putt was sunk, 0 otherwise
- x_i : distance to the hole of the i^{th} putt

Visualization: The first thing we should do is plot the data, which we do in Figure 4.1. What do you notice?

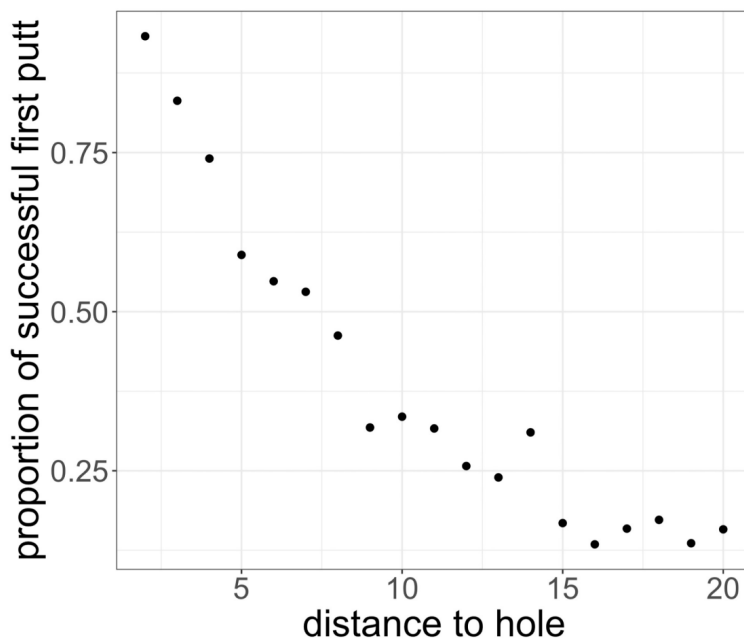


Figure 4.1: Scatter plot of the putt data

Like we would expect, the observed proportion of putts sunk decreases as the distance to the hole increases. Let's try to use what we've learned about linear regression to model this relationship.

4.1.2 Modeling

Model 1: Simple Linear Regression

We know from Lecture 1 that the this model takes the form:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \mathbb{E}[\epsilon_i] = 0 \quad (4.1)$$

We already know how to estimate the coefficients for this model. Recall that the least squares solution to this model is given by:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad (4.2)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad (4.3)$$

We fit this model in R with the following code:

```
model_1 = lm(putt_made ~ distance, data = putt_data)  
summary(model_1)
```

When we plot the fitted line against the data in Figure 4.2, we see that the model does not fit the data well.

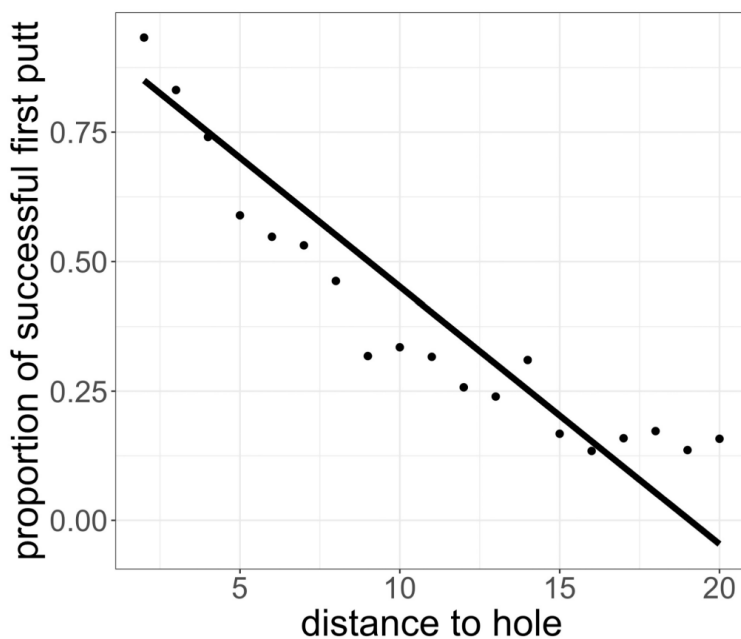


Figure 4.2: Putt data and fitted simple linear regression

We then move on to a more complex model: a cubic regression similar to the quadratic model we fit in Lecture 2.

Model 2: Cubic Regression

We know the form of this model as well:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \epsilon_i, \quad \mathbb{E}[\epsilon_i] = 0 \quad (4.4)$$

We also know how to estimate these coefficients from Lecture 2. Recall that the least squares solution to this model is given by:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (4.5)$$

We fit this model in R with the following code:

```
model_2 = lm(putt_made ~ distance + I(distance^2) + I(distance^3),
             data = putt_data)
summary(model_2)
```

When we plot the fitted line against the data in Figure 4.3, we see that the model fits the data much better when $x_i \in [0, 20]$.

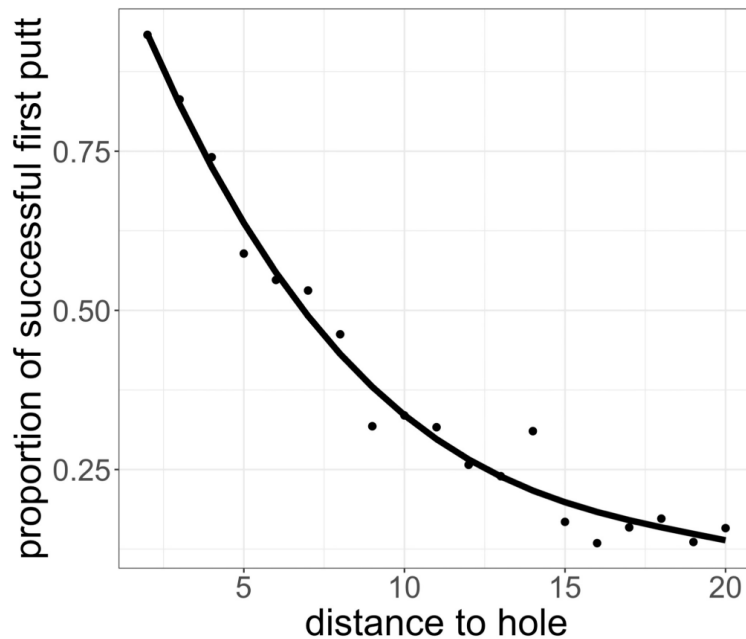


Figure 4.3: Putt data and fitted cubic regression

However, when we zoom out in Figure 4.4, we see that the model is unable to extrapolate when $x_i > 20$: in fact, the estimated probability of sinking the putt goes below 0.

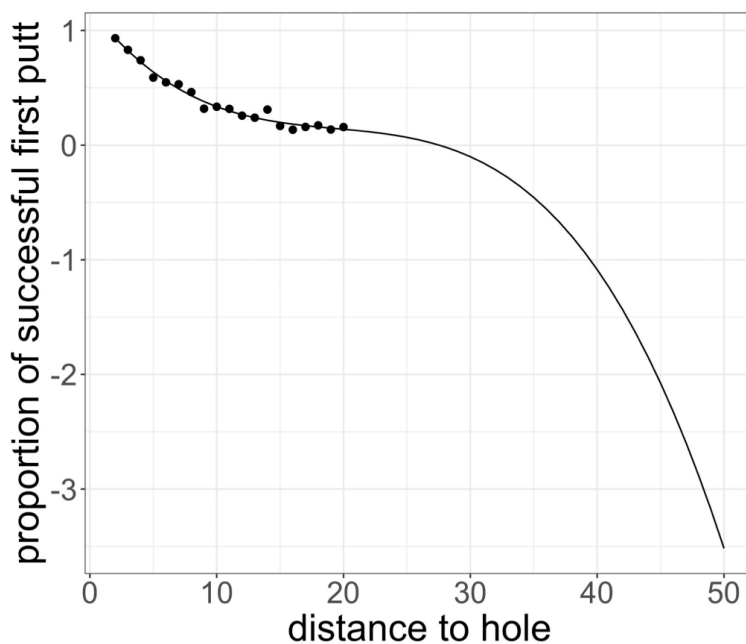


Figure 4.4: Putt data and fitted cubic regression (zoomed out)

This presents a problem: we need our predictions \hat{y}_i to be between 0 and 1, but ordinary linear regression does **not** guarantee this. How can we constrain our model to ensure that our predictions in $[0, 1]$?

4.2 Logistic Regression

4.2.1 The Logistic Function

We are looking for a "squishification" function that takes numbers in \mathbb{R} and maps them to numbers in $[0, 1]$. One such function is the **logistic (or sigmoid) function**:

Definition 4.1 (Logistic Function). *The **logistic function** is defined as:*

$$\text{Logistic}(z) = \text{sigmoid}(z) = \sigma(z) = \frac{1}{1 + e^{-z}} \quad (4.6)$$

We plot the logistic function in Figure 4.5.

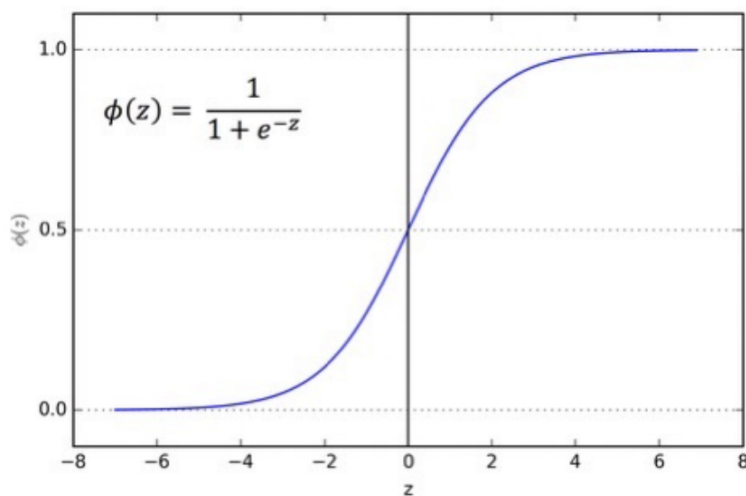


Figure 4.5: The Logistic Function

How can we implement this in our regression modeling?

4.2.2 The Model

Before, we had a model of the form:

$$\hat{y}_i = \beta_0 + \beta_1 x_i \quad (4.7)$$

Now, we modify this model so that our predictions \hat{y}_i are between 0 and 1. We can do this by making our response variable \hat{y}_i the output of the logistic function:

$$\hat{y}_i = \text{Logistic}(\beta_0 + \beta_1 x_i) = \frac{1}{1 + e^{-\beta_0 + \beta_1 x_i}} \quad (4.8)$$

This is a **logistic regression model**, and allows us to model the **probability** that a putt is sunk directly. Mathematically, we can write this model as:

$$p_i = \mathbb{P}(y_i = 1 \mid x_i) = \frac{1}{1 + e^{-(x_i^T \boldsymbol{\beta})}} \quad (4.9)$$

$$\text{where } y_i \sim \text{Bernoulli}(p_i) = \begin{cases} 1 & \text{w.p. } p_i \\ 0 & \text{w.p. } 1 - p_i \end{cases} \quad (4.10)$$

4.2.3 Estimating the Coefficients

Our data is in terms of $y_i \in \{0, 1\}$ and $x_i \in \mathbb{R}$, not in terms of $p_i \in [0, 1]$. How can we then estimate the coefficients $\boldsymbol{\beta}$ in logistic regression?

In linear regression, we estimate $\boldsymbol{\beta}$ by minimizing the Residual Sum of Squares (RSS), or the squared error.

Definition 4.2 (Residual Sum of Squares). The **Residual Sum of Squares** (*RSS*) is defined as:

$$\begin{aligned} RSS &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - x_i^T \beta)^2 \end{aligned} \quad (4.11)$$

In logistic regression, we estimate β by minimizing the **log loss**, or the cross-entropy loss.

Definition 4.3 (Log Loss). The **log loss** is defined as:

$$\begin{aligned} L(\beta) &= -\frac{1}{n} \sum_{i=1}^n (y_i \log p_i + (1 - y_i) \log(1 - p_i)) \\ \text{where } p_i &= \mathbb{P}(y_i = 1 \mid x_i, \beta) = \frac{1}{1 + e^{-(x_i^T \beta)}} \end{aligned} \quad (4.12)$$

1. If $y_i = 1$, then $L(\beta) = -\log p_i$
 - If $p_i \approx 1$, then $\log p_i$ is high and $L(\beta)$ is low
 - If $p_i \approx 0$, then $\log p_i$ is low and $L(\beta)$ is high
2. If $y_i = 0$, then $L(\beta) = -\log(1 - p_i)$
 - If $p_i \approx 1$, then $\log(1 - p_i)$ is low and $L(\beta)$ is high
 - If $p_i \approx 0$, then $\log(1 - p_i)$ is high and $L(\beta)$ is low

The log loss is a **convex** function, which means that it has a single global minimum. This means we can solve for β by setting its gradient equal to 0.

$$\begin{aligned} \nabla_{\beta} L(\beta) &= -\frac{1}{n} \sum_{i=1}^n (y_i \nabla_{\beta} \log p_i + (1 - y_i) \nabla_{\beta} \log(1 - p_i)) \\ &= 0 \end{aligned} \quad (4.13)$$

Solving for β analytically can be difficult, so we use gradient descent in R to find the minimum iteratively.

4.3 Back to Sinking Putts

Recall the variables we have:

- i : index of the i^{th} putt in our dataset
- y_i : 1 if the i^{th} putt was sunk, 0 otherwise
- x_i : distance to the hole of the i^{th} putt

4.3.1 Fitting the Model

We use the `glm` function to fit the logistic regression model.

```
model = glm(putt_made ~ distance, data = putt_data, family = "binomial")  
summary(model)
```

We plot the fitted line against the data in Figure 4.6. It seems to fit the data well, but we should check the model's performance from longer distances.

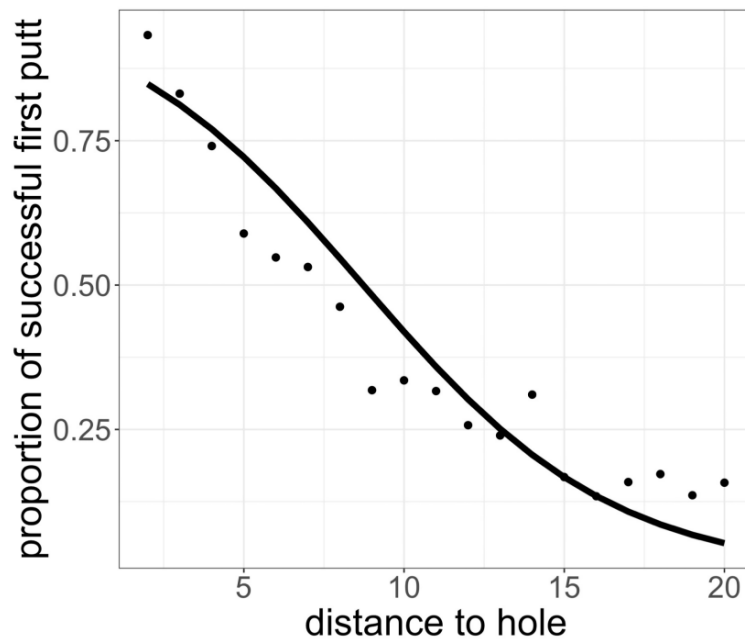


Figure 4.6: Putt data and fitted logistic regression

In Figure 4.7, we see that the model is able to extrapolate when $x_i > 20$ because we forced our output to be between 0 and 1.

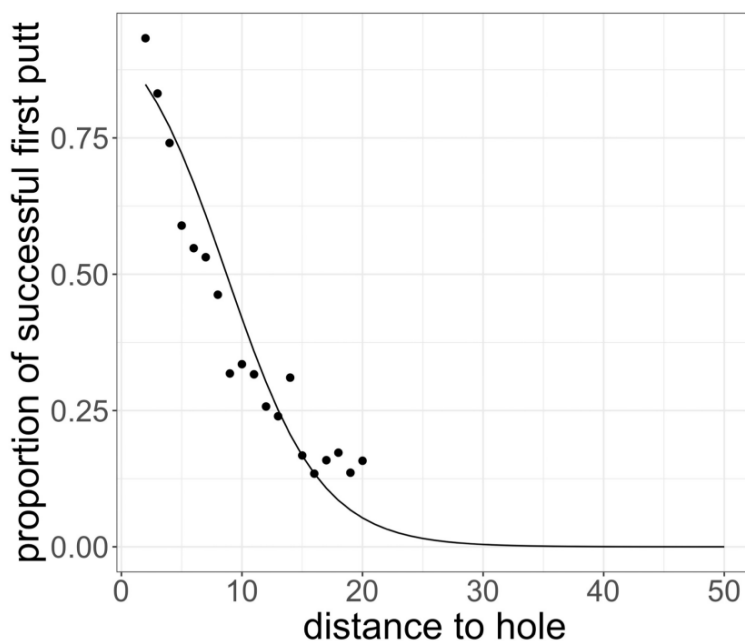


Figure 4.7: Putt data and fitted logistic regression (zoomed out)

4.3.2 Improving the Model

We can do even better by modeling the log odds with a cubic polynomial. Our new model is:

$$\mathbb{P}(y_i = 1 \mid x_i) = \frac{1}{1 + e^{-\beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3}} \quad (4.14)$$

We fit this model in R with the following code:

```
model = glm(putt_made ~ distance + I(distance^2) + I(distance^3),  
             data = putt_data, family = "binomial")  
summary(model)
```

We plot the fitted line against the data in Figure 4.8, and we can see that the model is able to extrapolate when $x_i > 20$.

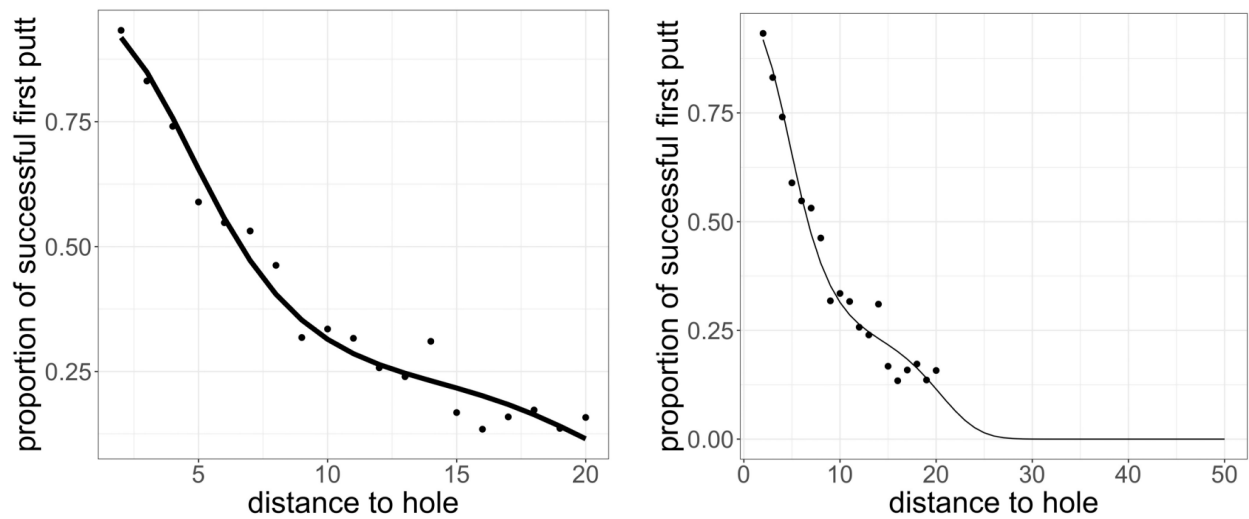


Figure 4.8: Putt data and fitted logistic regression with a cubic polynomial

4.3.3 Takeaways

What can we take away from this example?

- Use linear regression to predict a value that can take any real number
- Use logistic regression to predict a **probability** in $[0, 1]$