# Logistic Regression Introduction

## Julie Deeke
Statistics with Python Course Developer

# Cartwheel Data

Random sample of 25 adults attempted cartwheels
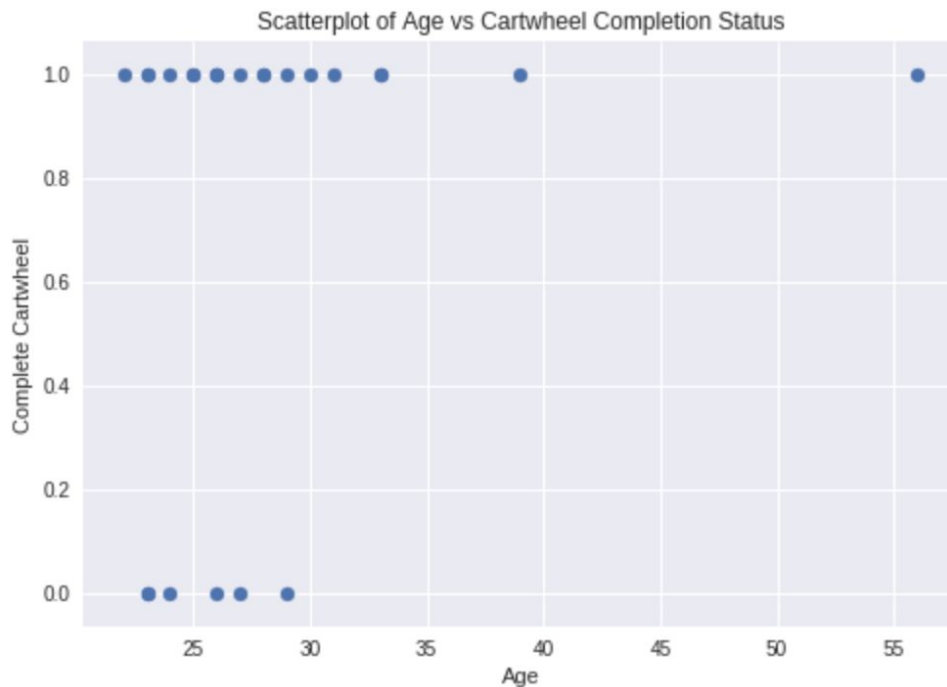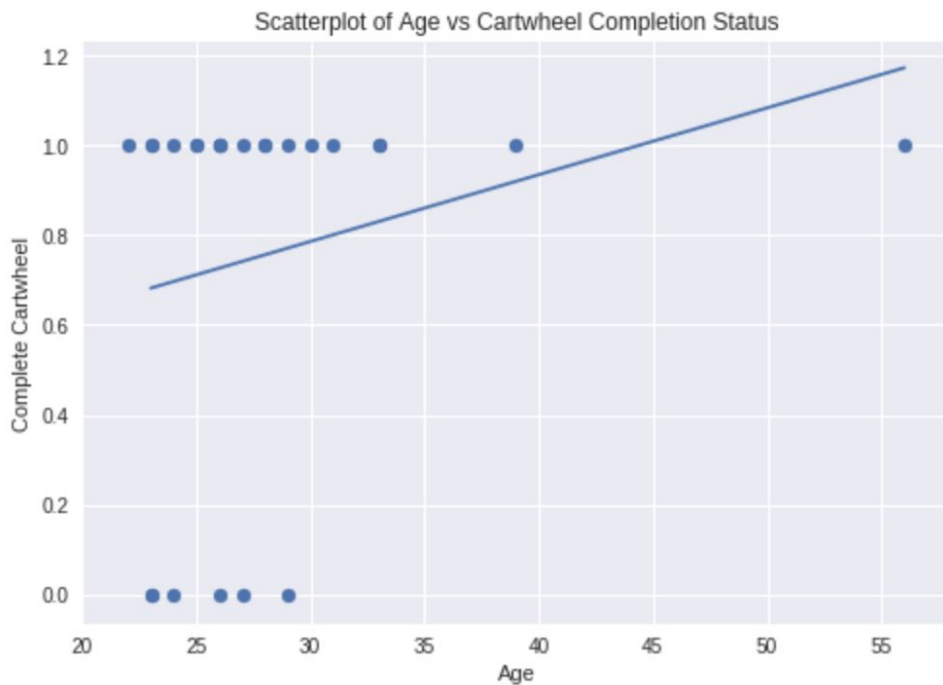
**Primary Variable of interest:** Cartwheel completion

# Research Question

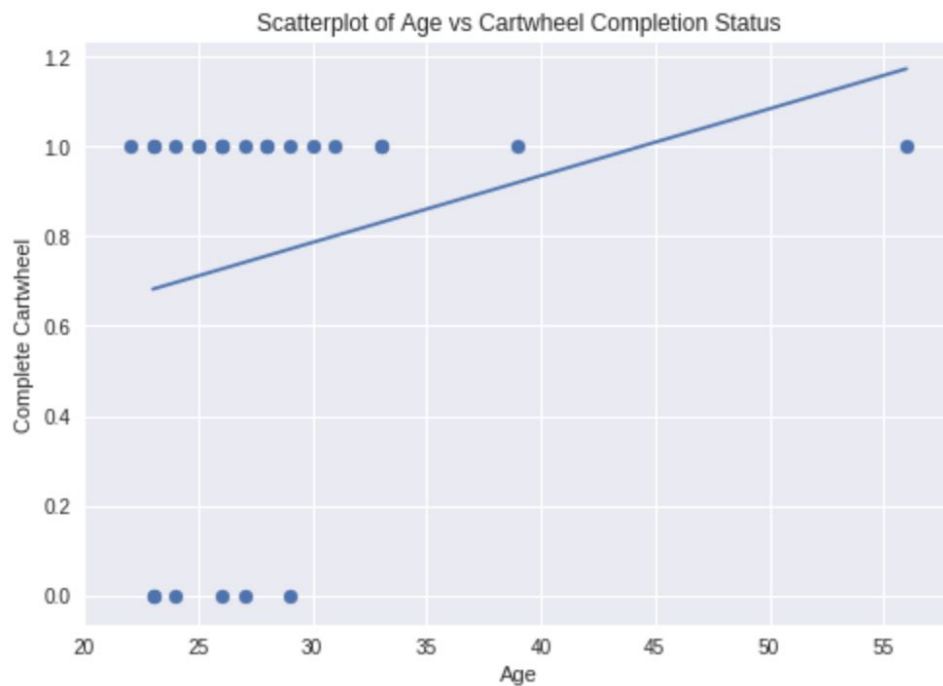Based on age, can we predict whether a cartwheel is completed?

# Let's Look at the Data

# Linear Model



Scatterplot of Age vs Cartwheel Completion Status

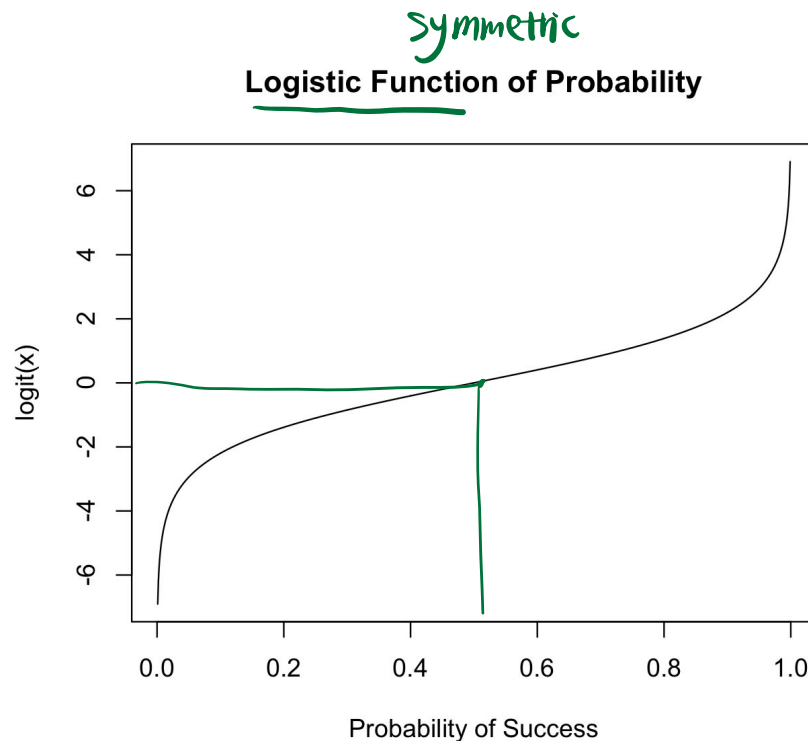# Linear Model



Scatterplot of Age vs Cartwheel Completion Status

$\hat{y} = 0.34 + 0.015$ age

# Logit Transformation

- Instead of predicting completion status, we predict a ***transformed version*** of the probability of a success

# Logit Transformation

- Instead of predicting completion status, we predict a **_transformed version_** of the probability of a success

- Uses the logit function $\ln\left(\dfrac{p}{1-p}\right)$

*Symmetric*

**Logistic Function of Probability**
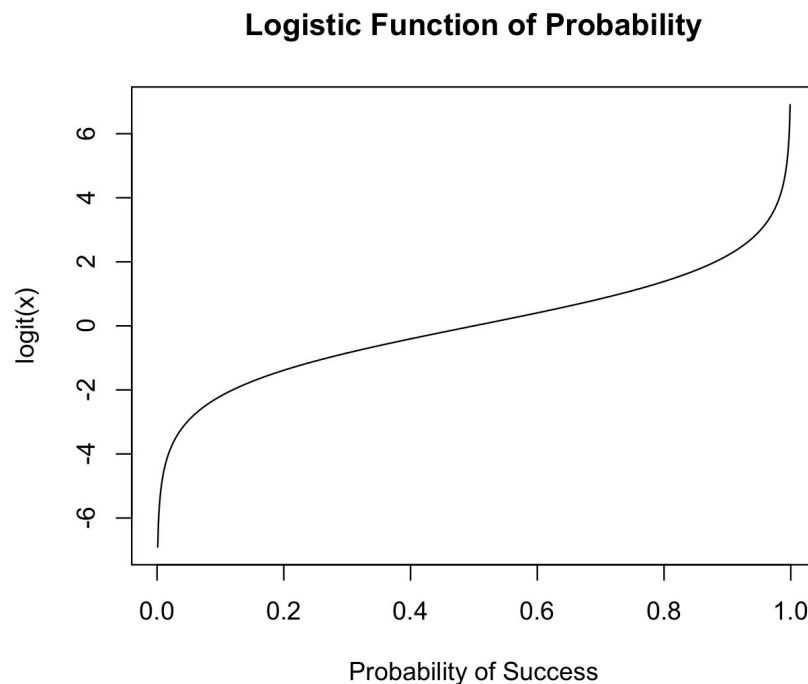


logit(x) vs Probability of Success

# Logit Transformation

- Instead of predicting completion status, we predict a ***transformed version*** of the probability of a success

- Uses the logit function:

$$\ln\left(\frac{p}{1-p}\right)$$

- $\mathrm{logit}(\hat{y}) = b_0 + b_1 x$

**Logistic Function of Probability**



y-axis: logit(x) — values -6, -4, -2, 0, 2, 4, 6

x-axis: Probability of Success — values 0.0, 0.2, 0.4, 0.6, 0.8, 1.0

# Logistic Regression Line



Scatterplot of Age vs. Cartwheel Completion Status

# Logistic Regression Line



Scatterplot of Age vs. Cartwheel Completion Status

# Extrapolation IVQ

**Would you feel comfortable using this model to estimate the probability that a teenager who is 15 can complete a cartwheel?**

# Logistic Regression Equation

Generalized Linear Model Regression Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | CompleteGroup | **No. Observations:** | 25 |
| **Model:** | GLM | **Df Residuals:** | 23 |
| **Model Family:** | Binomial | **Df Model:** | 1 |
| **Link Function:** | logit | **Scale:** | 1.0 |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **Intercept** | -4.4213 | 4.429 | -0.998 | 0.318 | -13.101 | 4.259 |
| **Age** | 0.2096 | 0.171 | 1.225 | 0.221 | -0.126 | 0.545 |

# Logistic Regression Equation

Generalized Linear Model Regression Results

| Dep. Variable: | CompleteGroup | No. Observations: | 25 |
|---|---|---|---|
| Model: | GLM | Df Residuals: | 23 |
| Model Family: | Binomial | Df Model: | 1 |
| Link Function: | logit | Scale: | 1.0 |

| | coef | std err | z | P>\|z\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -4.4213 | 4.429 | -0.998 | 0.318 | -13.101 | 4.259 |
| Age | 0.2096 | 0.171 | 1.225 | 0.221 | -0.126 | 0.545 |

# Logistic Regression Equation

logit( $\hat{y}$ ) = -4.42 + 0.2096 age

Generalized Linear Model Regression Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | CompleteGroup | **No. Observations:** | 25 |
| **Model:** | GLM | **Df Residuals:** | 23 |
| **Model Family:** | Binomial | **Df Model:** | 1 |
| **Link Function:** | logit | **Scale:** | 1.0 |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| **Intercept** | -4.4213 | 4.429 | -0.998 | 0.318 | -13.101 | 4.259 |
| **Age** | 0.2096 | 0.171 | 1.225 | 0.221 | -0.126 | 0.545 |

# Logistic Regression Equation

logit( $\hat{y}$ ) = -4.42 + 0.2096 age

**Slope interpretation:**
For each increase in age by 1 year, the log odds of a successful cartwheel increases by about 0.2096, on average.

Generalized Linear Model Regression Results

| Dep. Variable: | CompleteGroup | No. Observations: | 25 |
|---|---|---|---|
| Model: | GLM | Df Residuals: | 23 |
| Model Family: | Binomial | Df Model: | 1 |
| Link Function: | logit | Scale: | 1.0 |

|  | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -4.4213 | 4.429 | -0.998 | 0.318 | -13.101 | 4.259 |
| Age | 0.2096 | 0.171 | 1.225 | 0.221 | -0.126 | 0.545 |

# Logistic Regression Equation

logit( $\hat{y}$ ) = -4.42 + 0.2096 age

**Slope interpretation:** For each year increase in age, the odds of a successful cartwheel increases by about 1.23 ($e^{0.2096}$) times that of the younger age, on average.

Generalized Linear Model Regression Results

| | | | |
|---|---|---|---|
| **Dep. Variable:** | CompleteGroup | **No. Observations:** | 25 |
| **Model:** | GLM | **Df Residuals:** | 23 |
| **Model Family:** | Binomial | **Df Model:** | 1 |
| **Link Function:** | logit | **Scale:** | 1.0 |

| | coef | std err | z | P>|z| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -4.4213 | 4.429 | -0.998 | 0.318 | -13.101 | 4.259 |
| Age | 0.2096 | 0.171 | 1.225 | 0.221 | -0.126 | 0.545 |

# Predicted Probability of Success

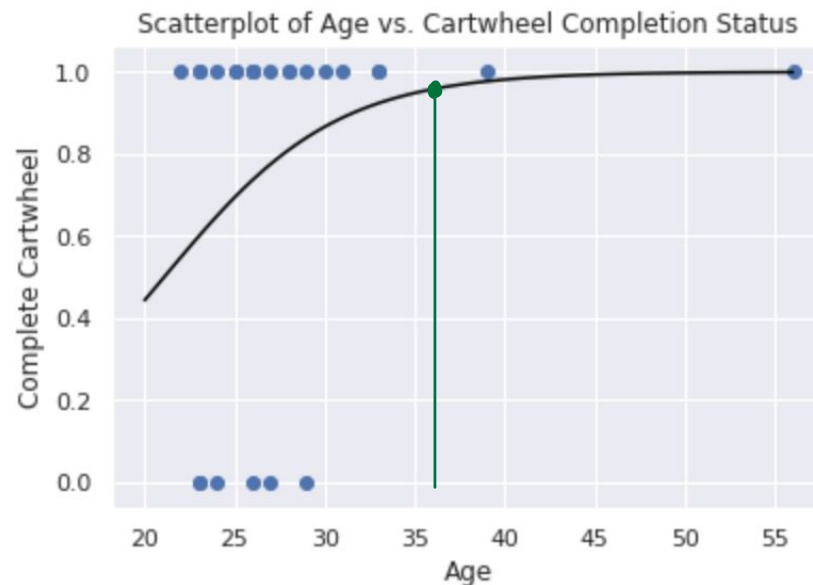- For someone who is 36, what is their predicted log odds of a successful cartwheel using the model?

# Predicted Probability of Success

- For someone who is 36, what is their predicted log odds of a successful cartwheel using the model?
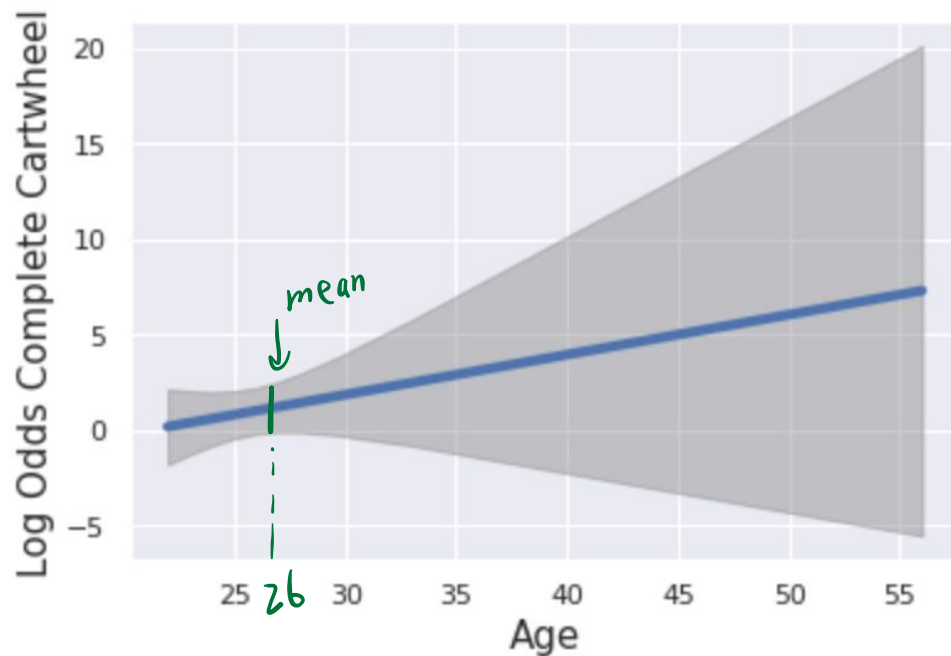
$$\text{logit}(\hat{y}) = -4.42 + 0.2096 \text{ age}$$
$$= -4.42 + 0.2096 \ (36)$$
$$= 3.13$$
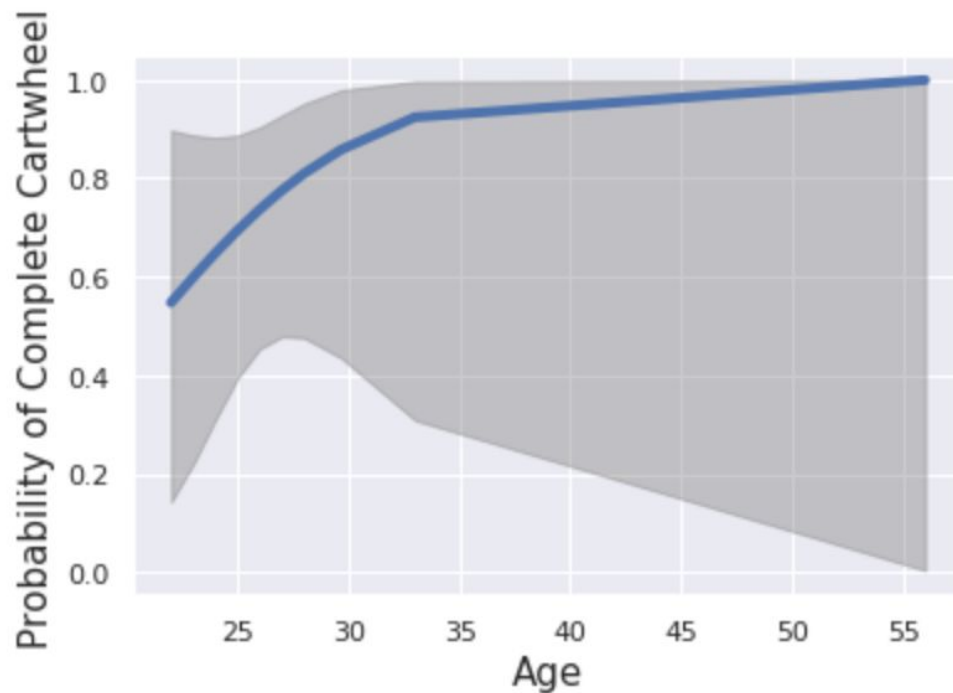
# Predicted Probability of Success

- For someone who is 36, what is their predicted log odds of a successful cartwheel using the model?
- Using the graph on the right, estimate what the probability of success might be?



Scatterplot of Age vs. Cartwheel Completion Status

# Prediction Uncertainty

# Prediction Uncertainty

# Assumptions

We need to assume that our model $\text{logit}(y) = \beta_0 + \beta_1 x_1$ is appropriate

# Assumptions

We need to assume that our model $\text{logit}(y) = \beta_0 + \beta_1 x_1$ is appropriate

~with a large enough sample size, you can identify discrepancies with residual plots

# Assumptions

We need to assume that our model $\text{logit}(y) = \beta_0 + \beta_1 x_1$ is appropriate

~with a large enough sample size, you can identify discrepancies with residual plots

~y only takes two values, so residuals can be limited

# Assumptions

We need to assume that our model $logit(y) = \beta_0 + \beta_1 x_1$ is appropriate

~with a large enough sample size, you can identify discrepancies with residual plots

~y only takes two values, so residuals can be limited

~to create informative residual plots, it helps if x takes a wide range of values and to have additional covariates