

Logistic Regression

Predicting whether a police stop results in a frisk

“The New York City stop-and-frisk policy allowed police to stop and search individuals when any officer had “reasonable suspicion” — a standard that is lower than the “probable cause” needed to justify an arrest. ...

Civil rights activists, such as the Center for Constitutional Rights, argued that this stop-and-frisk policy led to the New York City Police Department (NYPD) unfairly targeting people of color. Recent New York leaders such as Mayor Bill de Blasio made changes to that policy, decreasing the frequency of stop-and-frisks happening on a daily basis.”

<https://dataspace.sites.grinnell.edu/nypd1.html>

The response variable is binary: it takes two values

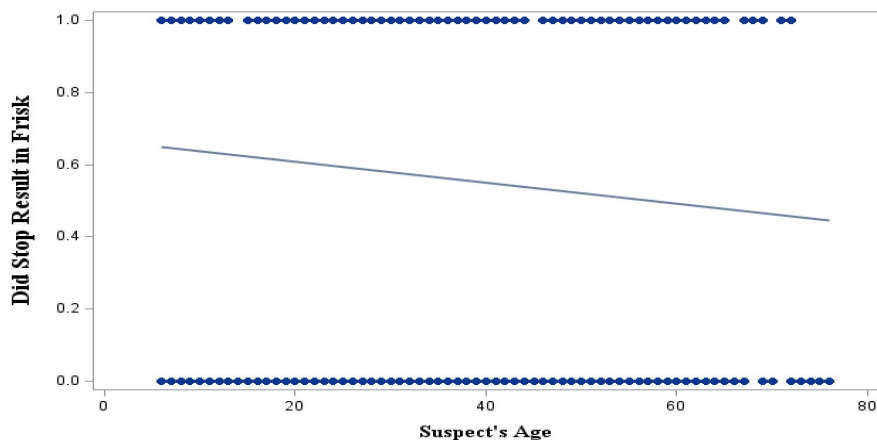
$Y = 1$ if the police stop results in a frisk

$Y = 0$ if the police stop does not result in a frisk

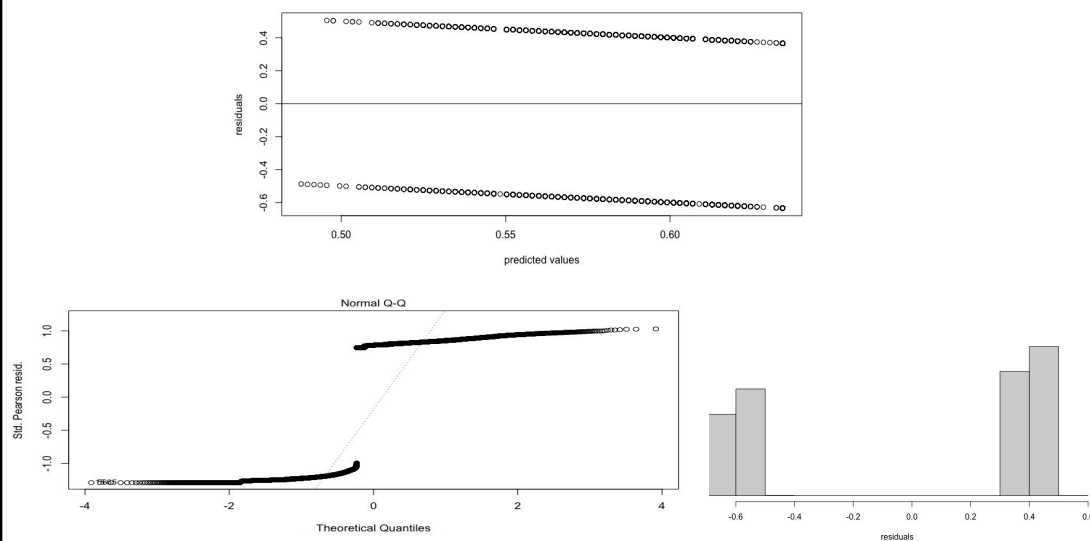
One explanatory variable, suspect’s age, is quantitative.

We will later look at the suspect’s sex and race as categorical explanatory variables.

What if we treat the *binary* response as *quantitative*?
i.e., what if we fit a simple linear regression model?



What if we treat the *binary* response as *quantitative*?
i.e., what if we fit a simple linear regression model?



Modeling the Mean Response

In linear regression, we model the mean response:

For a binary variable, what does the sample “mean” represent?

Status	
Frisk (Y = 1)	No Frisk (Y = 0)
6072	4122

- In logistic regression, we use the _____ to model the _____, which changes as a function of X .

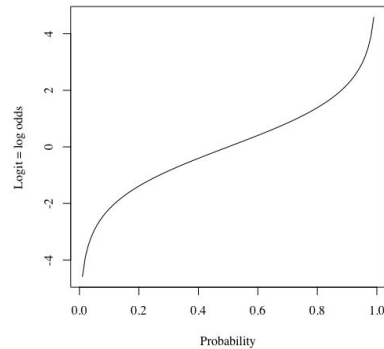
Using a transformation to fix problem

In logistic regression, instead of constructing a model for Y , we construct a model for _____

- Let π be the probability of “success”
- Since π is a probability, it can take values from _____.
- Possible model: _____
 - Problem: _____
- Second possible model: Instead of modeling the probability π directly, model the “odds” of success
- _____
 - Range of values: _____
 - Problem: _____

The Logit Transformation

The **logit function** is defined as



i.e., the logit is the _____.

Range of values of logit: _____

The Logistic Regression model is given by

The Logistic Regression Model

$$\text{Logit}(\pi) = \ln\left(\frac{\pi}{1 - \pi}\right) = \beta_0 + \beta_1 X$$

- The logit is a one-to-one transformation, i.e.,
 - For every value of π (with the exception of 0 and 1),
there is one, and only one, value of the logit
 - For every value of the logit (with the exception of $-\infty$ and $+\infty$),
there is one, and only one, value of π

Probability Form: We can rewrite our model so that the *probability of success* (π) is a function of the model parameters and the predictor X .

Four Probabilities (definitions on pg. 419)

Let $P(\text{success})$ = the probability of success. There are four version of this probability of interest:

	TRUE Value	FITTED Value
Actual Probability		
Model Probability		

If the model is exactly correct, then _____

Interpreting The Logistic Regression Parameters

In the logit form:

$$\ln\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 X$$

We are modeling the log-odds of “success”:

An age increase of 1 year for a suspect is associated with a log-odds of the police stop resulting in a frisk increase by _____ on average

We can calculate the “odds” of success by:

Interpretation

	coef	std err	z	P> z	[0.025	0.975]
Intercept	-0.3055	0.070	-4.346	0.000	-0.443	-0.168
arsenic	0.3791	0.039	9.840	0.000	0.304	0.455

Test and Confidence Interval for β_1

In the Simple Logistic Regression model,

The null and alternative hypotheses are given by:

The test statistic has the form:

The confidence interval for the β_1 :

Checking the Conditions

Conditions for Logistic Regression:

1. Linearity: The logit is a linear function of the explanatory variable X
2. Independence: The observations are not related to each other – no paired data, no clustered (grouped) data
3. Random: The data are obtained via a random process, either through (a) a random sample was obtained, or (b) random assignment to groups.

Conditions NOT required:

1. Constant Variance: Variance will NOT be constant!!
2. Normality: Population and sample will NOT be Normal...

Interpretation

	coef	std err	z	P> z	[0.025	0.975]
Intercept	0.0027	0.079	0.035	0.972	-0.153	0.158
arsenic	0.4608	0.041	11.134	0.000	0.380	0.542
dist100	-0.8966	0.104	-8.593	0.000	-1.101	-0.692