

Lecture 12: Linear Regression

Heidi Perry, PhD

Hack University

heidiperryphd@gmail.com

11/15/2016

Presentation derived from OpenIntro Statistics presentation for Chapter 7. These slides are available at <http://www.openintro.org> under a [Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported license \(CC BY-NC-SA\)](#).

- 1 Logistic Regression, A Generalized Linear Model
- 2 A Simple Example - The Donner Party
- 3 Interpretation of Coefficients - Odds Ratio

Odds

Odds are another way of quantifying the probability of an event, commonly used in gambling (and logistic regression).

Odds

For some event E ,

$$\text{odds}(E) = \frac{P(E)}{P(E^c)} = \frac{P(E)}{1 - P(E)}$$

Similarly, if we are told the odds of E are x to y then

$$\text{odds}(E) = \frac{x}{y} = \frac{x/(x+y)}{y/(x+y)}$$

which implies

$$P(E) = x/(x+y), \quad P(E^c) = y/(x+y)$$

Example - Donner Party

In 1846 the Donner and Reed families left Springfield, Illinois, for California by covered wagon. In July, the Donner Party, as it became known, reached Fort Bridger, Wyoming. There its leaders decided to attempt a new and untested route to the Sacramento Valley. Having reached its full size of 87 people and 20 wagons, the party was delayed by a difficult crossing of the Wasatch Range and again in the crossing of the desert west of the Great Salt Lake. The group became stranded in the eastern Sierra Nevada mountains when the region was hit by heavy snows in late October. By the time the last survivor was rescued on April 21, 1847, 40 of the 87 members had died from famine and exposure to extreme cold.

From Ramsey, F.L. and Schafer, D.W. (2002). *The Statistical Sleuth: A Course in Methods of Data Analysis* (2nd ed)

Example - Donner Party - Data

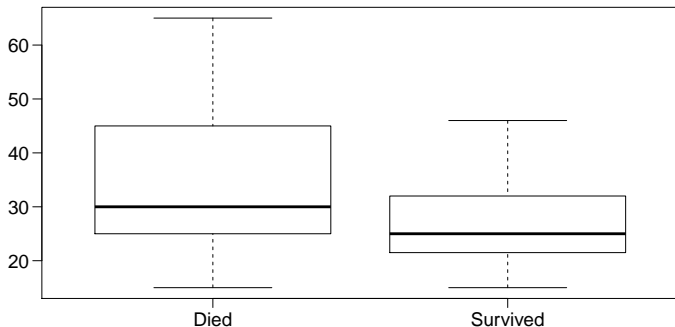
| | Age | Sex | Status |
|----|-------|--------|----------|
| 1 | 23.00 | Male | Died |
| 2 | 40.00 | Female | Survived |
| 3 | 40.00 | Male | Survived |
| 4 | 30.00 | Male | Died |
| 5 | 28.00 | Male | Died |
| ⋮ | ⋮ | ⋮ | ⋮ |
| 43 | 23.00 | Male | Survived |
| 44 | 24.00 | Male | Died |
| 45 | 25.00 | Female | Survived |

Example - Donner Party - EDA

Status vs. Gender:

| | Male | Female |
|----------|------|--------|
| Died | 20 | 5 |
| Survived | 10 | 10 |

Status vs. Age:



Example - Donner Party

It seems clear that both age and gender have an effect on someone's survival, how do we come up with a model that will let us explore this relationship?

Even if we set Died to 0 and Survived to 1, this isn't something we can transform our way out of - we need something more.

One way to think about the problem - we can treat Survived and Died as successes and failures arising from a binomial distribution where the probability of a success is given by a transformation of a linear model of the predictors.

Generalized linear models

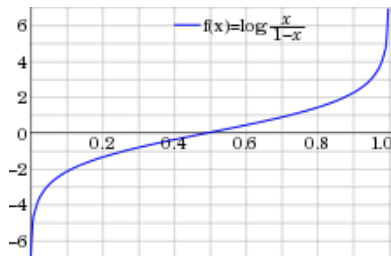
It turns out that this is a very general way of addressing this type of problem in regression, and the resulting models are called generalized linear models (GLMs). Logistic regression is just one example of this type of model.

All generalized linear models have the following three characteristics:

- 1 A probability distribution describing the outcome variable
- 2 A linear model
 - $\eta = \beta_0 + \beta_1 X_1 + \cdots + \beta_n X_n$
- 3 A link function that relates the linear model to the parameter of the outcome distribution
 - $g(p) = \eta$ or $p = g^{-1}(\eta)$

Logit Function

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right), \text{ for } 0 \leq p \leq 1$$

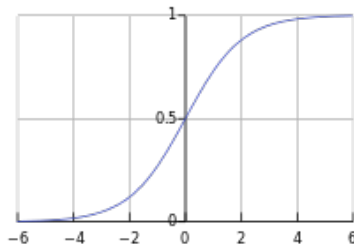


By Krishnavedala - Own work, CC0

Logistic Function

(a sigmoid curve)

$$g^{-1}(x) = \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{1 + \exp(-x)}$$



By Qef (talk) - Created from scratch with gnuplot,
Public Domain

The logistic regression model

The three GLM criteria give us:

$$y_i \sim \text{Binom}(p_i)$$

$$\eta = \beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n$$

$$\text{logit}(p) = \eta$$

From which we arrive at,

$$p_i = \frac{\exp(\beta_0 + \beta_1 x_{1,i} + \cdots + \beta_n x_{n,i})}{1 + \exp(\beta_0 + \beta_1 x_{1,i} + \cdots + \beta_n x_{n,i})}$$

Example - Donner Party - Model

See Jupyter notebook for code to run the regression.

Logit Regression Results

| | | | |
|-----------------------|------------------|--------------------------|---------|
| Dep. Variable: | Survived | No. Observations: | 45 |
| Model: | Logit | Df Residuals: | 43 |
| Method: | MLE | Df Model: | 1 |
| Date: | Wed, 16 Nov 2016 | Pseudo R-squ.: | 0.08954 |
| Time: | 20:16:29 | Log-Likelihood: | -28.145 |
| converged: | True | LL-Null: | -30.913 |
| | | LLR p-value: | 0.01863 |

| | coef | std err | z | P> z | [95.0% Conf. Int.] |
|------------------|---------|---------|--------|-------|--------------------|
| Age | -0.0665 | 0.032 | -2.063 | 0.039 | -0.130 -0.003 |
| intercept | 1.8185 | 0.999 | 1.820 | 0.069 | -0.140 3.777 |

Example - Donner Party - Prediction

| | Estimate | Std. Error | z value | Pr(> z) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 1.8185 | 0.9994 | 1.82 | 0.0688 |
| Age | -0.0665 | 0.0322 | -2.06 | 0.0391 |

Model:

$$\log \left(\frac{p}{1-p} \right) = 1.8185 - 0.0665 \times \text{Age}$$

Odds / Probability of survival for a newborn (Age=0):

$$\begin{aligned}\log \left(\frac{p}{1-p} \right) &= 1.8185 - 0.0665 \times 0 \\ \frac{p}{1-p} &= \exp(1.8185) = 6.16 \\ p &= 6.16 / 7.16 = 0.86\end{aligned}$$

Example - Donner Party - Prediction (cont.)

Model:

$$\log\left(\frac{p}{1-p}\right) = 1.8185 - 0.0665 \times \text{Age}$$

Odds / Probability of survival for a 25 year old:

$$\log\left(\frac{p}{1-p}\right) = 1.8185 - 0.0665 \times 25$$

$$\frac{p}{1-p} = \exp(0.156) = 1.17$$

$$p = 1.17/2.17 = 0.539$$

Odds / Probability of survival for a 50 year old:

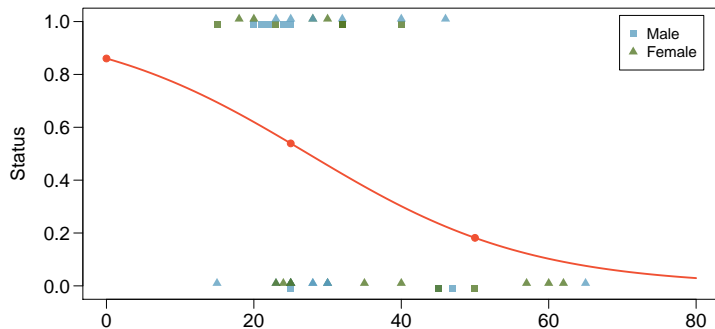
$$\log\left(\frac{p}{1-p}\right) = 1.8185 - 0.0665 \times 50$$

$$\frac{p}{1-p} = \exp(-1.5065) = 0.222$$

$$p = 0.222/1.222 = 0.181$$

Example - Donner Party - Prediction (cont.)

$$\log\left(\frac{p}{1-p}\right) = 1.8185 - 0.0665 \times \text{Age}$$



Example - Donner Party - Interpretation

| | Estimate | Std. Error | z value | $\Pr(> z)$ |
|-------------|----------|------------|---------|-------------|
| (Intercept) | 1.8185 | 0.9994 | 1.82 | 0.0688 |
| Age | -0.0665 | 0.0322 | -2.06 | 0.0391 |

Simple interpretation is only possible in terms of log odds and log odds ratios for intercept and slope terms.

Intercept: The log odds of survival for a party member with an age of 0. From this we can calculate the odds or probability, but additional calculations are necessary.

Slope: For a unit increase in age (being 1 year older) how much will the log odds ratio change, not particularly intuitive. More often than not we care only about sign and relative magnitude.

Example - Donner Party - Interpretation - Slope

$$\begin{aligned}\log\left(\frac{p_1}{1-p_1}\right) &= 1.8185 - 0.0665(x+1) \\ &= 1.8185 - 0.0665x - 0.0665\end{aligned}$$

$$\log\left(\frac{p_2}{1-p_2}\right) = 1.8185 - 0.0665x$$

$$\log\left(\frac{p_1}{1-p_1}\right) - \log\left(\frac{p_2}{1-p_2}\right) = -0.0665$$

$$\log\left(\frac{p_1}{1-p_1} \bigg/ \frac{p_2}{1-p_2}\right) = -0.0665$$

$$\frac{p_1}{1-p_1} \bigg/ \frac{p_2}{1-p_2} = \exp(-0.0665) = 0.94$$

Example - Donner Party - Age and Gender

Logit Regression Results

| | | | |
|-----------------------|------------------|--------------------------|----------|
| Dep. Variable: | Survived | No. Observations: | 45 |
| Model: | Logit | Df Residuals: | 42 |
| Method: | MLE | Df Model: | 2 |
| Date: | Wed, 16 Nov 2016 | Pseudo R-squ.: | 0.1710 |
| Time: | 20:27:45 | Log-Likelihood: | -25.628 |
| converged: | True | LL-Null: | -30.913 |
| | | LLR p-value: | 0.005066 |

| | coef | std err | z | P> z | [95.0% Conf. Int.] |
|------------------|---------|---------|--------|-------|--------------------|
| Age | -0.0782 | 0.037 | -2.097 | 0.036 | -0.151 -0.005 |
| Female | 1.5973 | 0.756 | 2.114 | 0.034 | 0.117 3.078 |
| intercept | 1.6331 | 1.110 | 1.471 | 0.141 | -0.543 3.809 |

Gender slope: When the other predictors are held constant this is the log odds ratio between the given level (Female) and the reference level (Male).

Example - Donner Party - Gender Models

Just like MLR we can plug in gender to arrive at two status vs age models for men and women respectively.

General model:

$$\log\left(\frac{p_1}{1-p_1}\right) = 1.63312 + -0.07820 \times \text{Age} + 1.59729 \times \text{Sex}$$

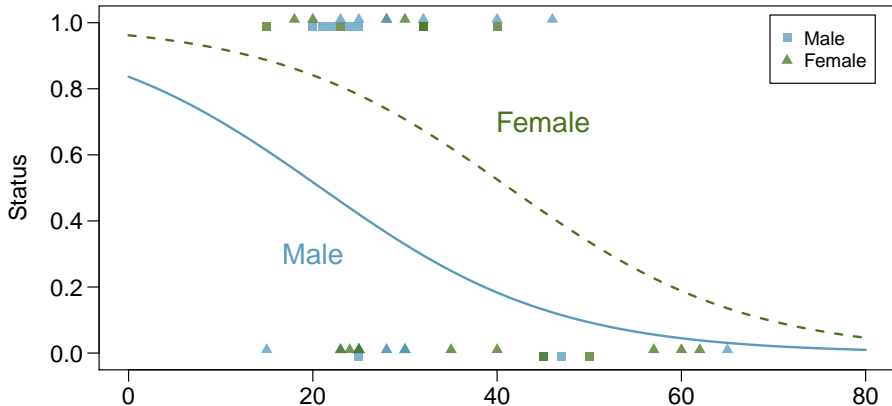
Male model:

$$\begin{aligned}\log\left(\frac{p_1}{1-p_1}\right) &= 1.63312 + -0.07820 \times \text{Age} + 1.59729 \times \mathbf{0} \\ &= 1.63312 + -0.07820 \times \text{Age}\end{aligned}$$

Female model:

$$\begin{aligned}\log\left(\frac{p_1}{1-p_1}\right) &= 1.63312 + -0.07820 \times \text{Age} + 1.59729 \times \mathbf{1} \\ &= 3.23041 + -0.07820 \times \text{Age}\end{aligned}$$

Example - Donner Party - Gender Models (cont.)



Confidence interval for age slope coefficient

| | Estimate | Std. Error | z value | Pr(> z) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 1.6331 | 1.1102 | 1.47 | 0.1413 |
| Age | -0.0782 | 0.0373 | -2.10 | 0.0359 |
| SexFemale | 1.5973 | 0.7555 | 2.11 | 0.0345 |

Remember, the interpretation for a slope is the change in log odds ratio per unit change in the predictor.

Log odds ratio:

$$CI = PE \pm CV \times SE = -0.0782 \pm 1.96 \times 0.0373 = (-0.1513, -0.0051)$$

Odds ratio:

$$\exp(CI) = (\exp -0.1513, \exp -0.0051) = (0.85960.9949)$$

Example - Birdkeeping and Lung Cancer

A 1972 - 1981 health survey in The Hague, Netherlands, discovered an association between keeping pet birds and increased risk of lung cancer. To investigate birdkeeping as a risk factor, researchers conducted a case-control study of patients in 1985 at four hospitals in The Hague (population 450,000). They identified 49 cases of lung cancer among the patients who were registered with a general practice, who were age 65 or younger and who had resided in the city since 1965. They also selected 98 controls from a population of residents having the same general age structure.

From Ramsey, F.L. and Schafer, D.W. (2002). *The Statistical Sleuth: A Course in Methods of Data Analysis* (2nd ed)

Example - Birdkeeping and Lung Cancer - Data

| | LC | FM | SS | BK | AG | YR | CD |
|-----|------------|--------|------|--------|-------|-------|-------|
| 1 | LungCancer | Male | Low | Bird | 37.00 | 19.00 | 12.00 |
| 2 | LungCancer | Male | Low | Bird | 41.00 | 22.00 | 15.00 |
| 3 | LungCancer | Male | High | NoBird | 43.00 | 19.00 | 15.00 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| 147 | NoCancer | Female | Low | NoBird | 65.00 | 7.00 | 2.00 |

LC Whether subject has lung cancer

FM Sex of subject

SS Socioeconomic status

BK Indicator for birdkeeping

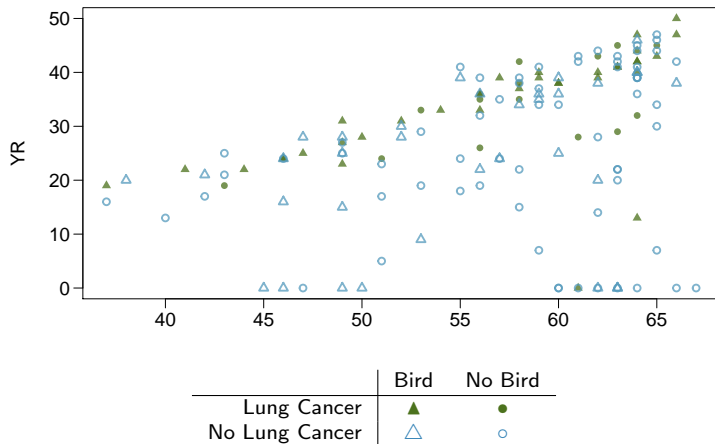
AG Age of subject (years)

YR Years of smoking prior to diagnosis or examination

CD Average rate of smoking (cigarettes per day)

NoCancer is the reference response (0 or failure), LungCancer is the non-reference response (1 or success) - this matters for interpretation.

Example - Birdkeeping and Lung Cancer - EDA



Example - Birdkeeping and Lung Cancer - Interpretation

| | Estimate | Std. Error | z value | Pr(> z) |
|-------------|----------|------------|---------|----------|
| (Intercept) | -1.9374 | 1.8043 | -1.07 | 0.2829 |
| FMFemale | 0.5613 | 0.5312 | 1.06 | 0.2907 |
| SSHigh | 0.1054 | 0.4688 | 0.22 | 0.8221 |
| BKBird | 1.3626 | 0.4113 | 3.31 | 0.0009 |
| AG | -0.0398 | 0.0355 | -1.12 | 0.2625 |
| YR | 0.0729 | 0.0265 | 2.75 | 0.0059 |
| CD | 0.0260 | 0.0255 | 1.02 | 0.3081 |

Keeping all other predictors constant then,

- The odds ratio of getting lung cancer for bird keepers vs non-bird keepers is $\exp(1.3626) = 3.91$.
- The odds ratio of getting lung cancer for an additional year of smoking is $\exp(0.0729) = 1.08$.

What do the numbers not mean ...

The most common mistake made when interpreting logistic regression is to treat an odds ratio as a ratio of probabilities.

Bird keepers are *not* 4x more likely to develop lung cancer than non-bird keepers.

This is the difference between relative risk and an odds ratio.

$$RR = \frac{P(\text{disease}|\text{exposed})}{P(\text{disease}|\text{unexposed})}$$

$$OR = \frac{P(\text{disease}|\text{exposed})/[1 - P(\text{disease}|\text{exposed})]}{P(\text{disease}|\text{unexposed})/[1 - P(\text{disease}|\text{unexposed})]}$$

Back to the birds

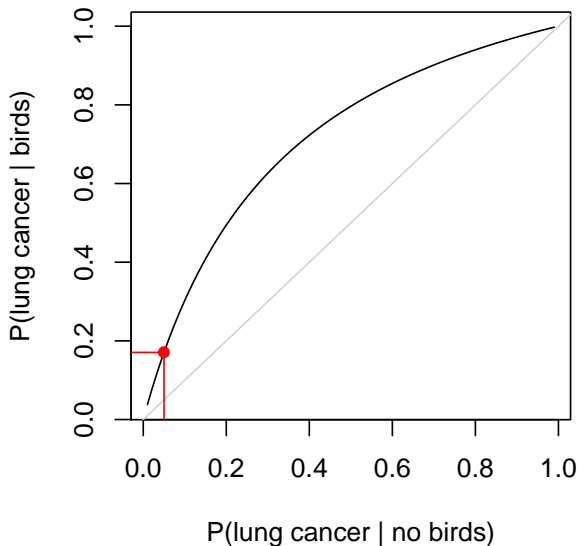
What is probability of lung cancer in a bird keeper if we knew that $P(\text{lung cancer}|\text{no birds}) = 0.05$?

$$\begin{aligned} OR &= \frac{P(\text{lung cancer}|\text{birds})/[1 - P(\text{lung cancer}|\text{birds})]}{P(\text{lung cancer}|\text{no birds})/[1 - P(\text{lung cancer}|\text{no birds})]} \\ &= \frac{P(\text{lung cancer}|\text{birds})/[1 - P(\text{lung cancer}|\text{birds})]}{0.05/[1 - 0.05]} = 3.91 \end{aligned}$$

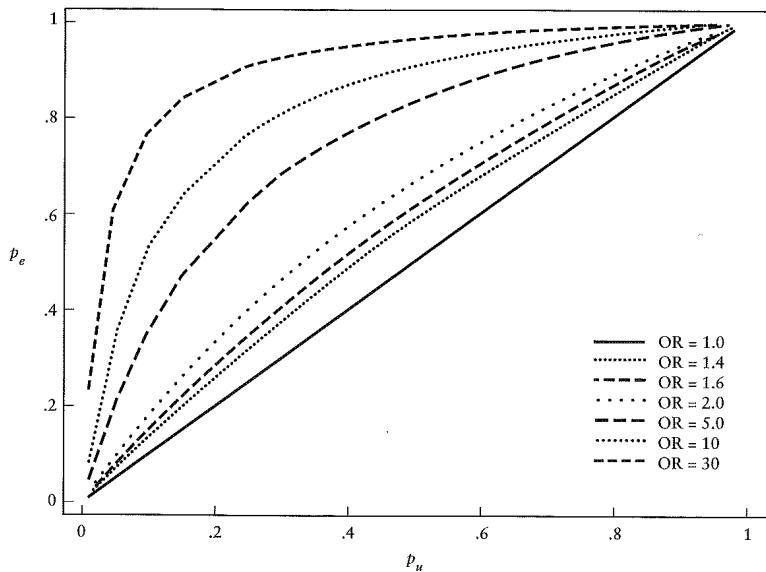
$$P(\text{lung cancer}|\text{birds}) = \frac{3.91 \times \frac{0.05}{0.95}}{1 + 3.91 \times \frac{0.05}{0.95}} = 0.171$$

$$RR = P(\text{lung cancer}|\text{birds})/P(\text{lung cancer}|\text{no birds}) = 0.171/0.05 = 3.41$$

Bird OR Curve



OR Curves





David Diez, Christopher Barr, & Mine Çetinkaya-Rundel (2015)

OpenIntro Statistics, [OpenIntro](#)

Recommended Reading

OpenIntro Statistics, Chapter 8

Data Science from Scratch, Chapter 16

Articles for discussion:

[Logistic Regression in Python Using Rodeo](#)

Example notebook: Lesson12_LogisticRegression.ipynb has code to generate models presented here.

No class next week.

Project updates when we return:

- 5 minutes
- Show data visualization and explain what you learned from it.
- Are you still working on your proposed data question, or after the exploratory analysis do you have a new question?
- Q&A/discussion: answer questions from the class (and ask questions!)