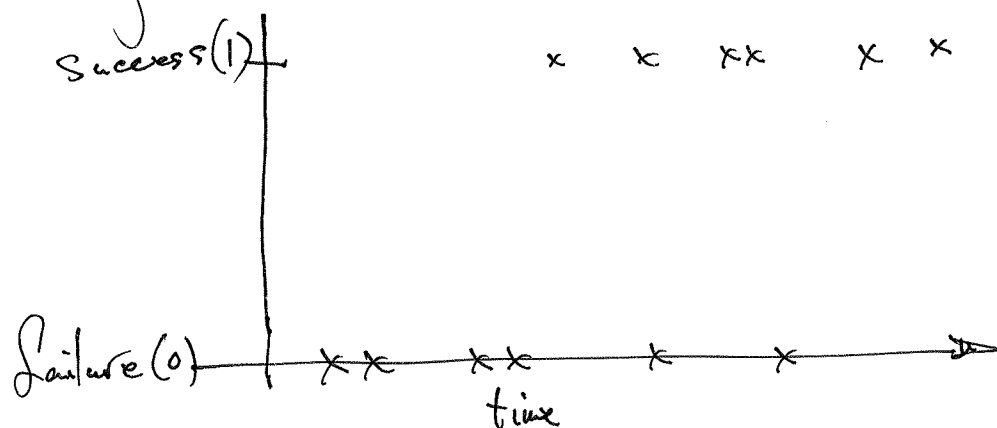


LOGISTIC REGRESSION

1

Suppose we want to predict the chance $p(x)$ of accomplishing a goal (like passing a course, or finishing a race in a certain amount of time) as a function of some quantitative variable x (like amount of study time, or time spent training). The response is binary, and data might look like this:



We seek a function $p(x)$ that sort-of "best fits" these data so that

Y is the response and has a Bernoulli(p) distribution (or could be binomial...)

$$p(x) = P(Y=1 | X=x)$$

(or polynomial)

Simple linear regression is not a good model in this case since we need the function $p(x)$ to satisfy $\lim_{x \rightarrow \infty} p(x) = 1$,

$\lim_{x \rightarrow -\infty} p(x) = 0$, and $p(x)$ should be monotone (note these are properties of cumulative distribution functions, so $p(x)$ is a cdf.)

2

Consider the odds of an event E :

$$\text{odds of } E := \frac{P(E)}{1-P(E)} \in [0, \infty].$$

While probabilities are always in $[0, 1]$, odds can take values from the extended non-negative real numbers...

The odds of an event CAN BE ∞ ... Now, the natural log of the odds ~~obey~~ obeys

$$-\infty \leq \log \frac{P(E)}{1-P(E)} \leq \infty$$

So maybe it would make sense to build a linear model for the natural log of the odds ratio:

$$\log \left(\frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x$$

After some algebra...

$$\frac{p(x)}{1-p(x)} = e^{\beta_0 + \beta_1 x}$$

$$\Rightarrow p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

or

$$p(x) = \left(1 + e^{-\beta_0 - \beta_1 x} \right)^{-1}$$

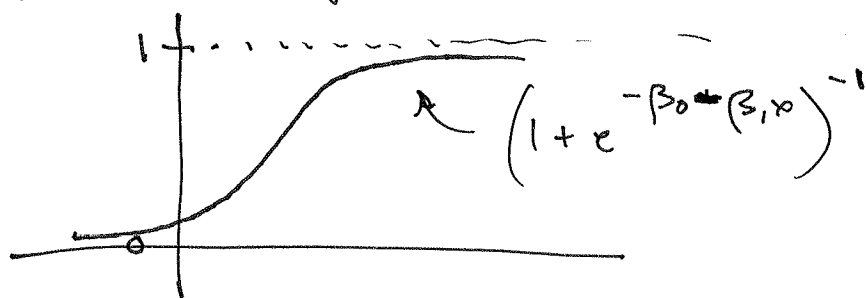


You should graph several instances of this $p(x)$ function (3) with different values of β_0 and β_1 , like

β_0	0	0	0	0	2	2	2	2	4	4	4	4
β_1	1	2	3	-4	1	2	3	-4	1	2	3	-4

When you do, you'll see that β_0 is a location parameter, and β_1 is a shape parameter. What happens as β_0 increases? as β_1 increases? when $\beta_1 < 0$?

Regardless of the β_0 and β_1 values, these curves have the same basic shape, which is a sigmoid or s-shape:



Like the cdf. of a normal random variable, this curve has rotational symmetry. Specifically, the point of rotational symmetry here is $(-\beta_0/\beta_1, 1/2)$.

There are other ways to transform $p(x)$ for which a linear model might be appropriate. We are mainly concerned with these three ways (they seem to be the most popular): (4)

① The logit function:

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

(which is what we've just been using)

② The probit function:

$$\eta(x) = \Phi^{-1}(x)$$

where $\Phi(x) = P(\text{a standard normal r.v.} \leq x)$.

③ The complementary log-log function

$$\eta(x) = \log(-\log(1-x))$$

The transformation η is typically called a link function (as in, the logit-link, or the probit link, etc.)

Back to the model in 😊 that uses the logit link...

How do we estimate the β_0 and β_1 coefficients? We typically need a computer to carry out a numerical procedure,

but to shed some light... we could write

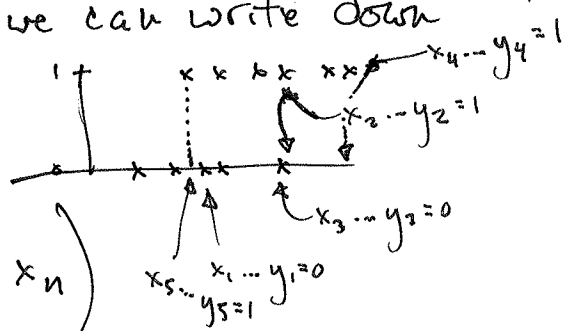
(5)

$$P(Y=y|X=x) = p(x)^y (1-p(x))^{1-y} \text{ for } y=0 \text{ or } 1.$$

Y is the Bernoulli($p(x)$)

response: $y \in \{0, 1\}$

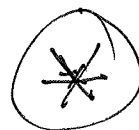
Now suppose we make n independent observations of Y when $X = x_1, x_2, \dots, x_n$. Then we can write down the likelihood function as



$$P(Y_1=y_1, \dots, Y_n=y_n | X_1=x_1, \dots, X_n=x_n)$$

by independence

$$= \prod_{k=1}^n P(Y_k=y_k | X_k=x_k)$$
$$= \prod_{k=1}^n p(x_k)^{y_k} (1-p(x_k))^{1-y_k}$$



We want to maximize the chance of seeing what we actually saw over all model parameter values (β_0 and β_1). This is called maximum likelihood. So fix the x_k and y_k

(6)

values in the conditional probability / likelihood function in $(*)$, and consider that function to be a function of the β_0 and β_1 . ~~We basically want to differentiate the~~ It turns out (as the case often is) it's easier to maximize the log of the likelihood function (which is okay since the log function is monotone). We call this the log-likelihood function:

$$L = \log P(Y_1 = y_1, \dots, Y_n = y_n | X_1 = x_1, \dots, X_n = x_n, \beta_0, \beta_1)$$

We need to think of the log-likelihood function as a function of the model parameters --- the data values are fixed.

$$= \log \prod_{k=1}^n [p(x_k)^{y_k} (1-p(x_k))^{1-y_k}] = \sum_{k=1}^n [\log p(x_k)^{y_k} + \log (1-p(x_k))^{1-y_k}]$$

$$= \sum_{k=1}^n y_k \log p(x_k) + \sum_{k=1}^n (1-y_k) \log (1-p(x_k))$$

$$= \sum_{k=1}^n \log (1-p(x_k)) + \sum_{k=1}^n y_k (\log (p(x_k)) - \log (1-p(x_k)))$$

$$= \sum_{k=1}^n \underbrace{\log (1-p(x_k))}_{= -\log [1 + e^{\beta_0 + \beta_1 x_k}]} + \sum_{k=1}^n y_k \underbrace{\log \frac{p(x_k)}{1-p(x_k)}}_{= \beta_0 + \beta_1 x_k}$$

7

And so our problem becomes finding the values of $\hat{\beta}_{0,MLE}$ and $\hat{\beta}_{1,MLE}$ that maximize

"maximum likelihood estimator"

$$\sum_{k=1}^n y_k \left(\hat{\beta}_{0,MLE} + \hat{\beta}_{1,MLE} x_k \right) - \sum_{k=1}^n \log \left[1 + e^{\hat{\beta}_{0,MLE} + \hat{\beta}_{1,MLE} x_k} \right].$$

This is difficult to do by hand! You can try ... try differentiating this expression with respect to $\hat{\beta}_{0,MLE}$ and $\hat{\beta}_{1,MLE}$ and setting to 0...

Statistical software packages (like SAS and R) will do this.