

Logistic regression

In statistics, the **logistic model** (or **logit model**) is a widely used statistical model that in its basic form uses a logistic function to model a binary dependent variable, although many more complex extensions exist. In regression analysis, **logistic regression** (or **logit regression**) is estimating the parameters of a logistic model (a form of binary regression). Mathematically, a binary logistic model has a dependent variable with two possible values, such as pass/fail, win/lose, alive/dead or healthy/sick; these are represented by an indicator variable, where the two values are labeled "0" and "1". In the logistic model, the log-odds (the logarithm of the odds) for the value labeled "1" is a linear combination of one or more independent variables ("predictors"); the independent variables can each be a binary variable (two classes, coded by an indicator variable) or a continuous variable (any real value). The corresponding probability of the value labeled "1" can vary between 0 (certainly the value "0") and 1 (certainly the value "1"), hence the labeling; the function that converts log-odds to probability is the logistic function, hence the name. The unit of measurement for the log-odds scale is called a *logit*, from ***log**istic **unit***, hence the alternative names. Analogous models with a different sigmoid function instead of the logistic function can also be used, such as the probit model; the defining characteristic of the logistic model is that increasing one of the independent variables multiplicatively scales the odds of the given outcome at a *constant* rate, with each dependent variable having its own parameter; for a binary independent variable this generalizes the odds ratio.

The binary logistic regression model has extensions to more than two levels of the dependent variable: categorical outputs with more than two values are modeled by multinomial logistic regression, and if the multiple categories are ordered, by ordinal logistic regression, for example the proportional odds ordinal logistic model.^[1] The model itself simply models probability of output in terms of input, and does not perform statistical classification (it is not a classifier), though it can be used to make a classifier, for instance by choosing a cutoff value and classifying inputs with probability greater than the cutoff as one class, below the cutoff as the other; this is a common way to make a binary classifier. The coefficients are generally not computed by a closed-form expression, unlike linear least squares; see § Model fitting. The logistic regression as a general statistical model was originally developed and popularized primarily by Joseph Berkson,^[2] beginning in Berkson (1944), where he coined "logit"; see § History.

Contents

Applications
<div> <div></div> <div>Examples</div> <div> <div>Logistic model</div> <div>Probability of passing an exam versus hours of study</div> </div> </div>
Discussion
<div> <div></div> <div>Logistic regression vs. other approaches</div> </div>
<div> <div></div> <div>Latent variable interpretation</div> </div>
<div> <div></div> <div>Logistic function, odds, odds ratio, and logit</div> <div> <div>Definition of the logistic function</div> <div>Definition of the inverse of the logistic function</div> <div>Interpretation of these terms</div> <div>Definition of the odds</div> <div>The odds ratio</div> <div>Multiple explanatory variables</div> </div> </div>
<div> <div></div> <div>Model fitting</div> <div> <div>"Rule of ten"</div> <div>Maximum likelihood estimation</div> <div>Iteratively reweighted least squares (IRLS)</div> <div>Evaluating goodness of fit</div> </div> </div>
<div> <div></div> <div>Coefficients</div> <div> <div>Likelihood ratio test</div> <div>Wald statistic</div> <div>Case-control sampling</div> </div> </div>
<div> <div></div> <div>Formal mathematical specification</div> <div> <div>Setup</div> <div>As a generalized linear model</div> <div>As a latent-variable model</div> <div>Two-way latent-variable model</div> <div>As a "log-linear" model</div> <div>As a single-layer perceptron</div> <div>In terms of binomial data</div> </div> </div>
<div> <div></div> <div>Bayesian</div> </div>
<div> <div></div> <div>History</div> </div>
<div> <div></div> <div>Extensions</div> </div>
<div> <div></div> <div>Software</div> </div>
<div> <div></div> <div>See also</div> </div>
<div> <div></div> <div>References</div> </div>
<div> <div></div> <div>Further reading</div> </div>
<div> <div></div> <div>External links</div> </div>

Applications

Logistic regression is used in various fields, including machine learning, most medical fields, and social sciences. For example, the Trauma and Injury Severity Score (TRISS), which is widely used to predict mortality in injured patients, was originally developed by Boyd et al. using logistic regression.^[3] Many other medical scales used to assess severity of a patient have been developed using logistic regression.^{[4][5][6][7]} Logistic regression may be used to predict the risk of developing a given disease (e.g. diabetes; coronary heart disease), based on observed characteristics of the patient (age, sex, body mass index, results of various blood tests, etc.).^{[8][9]} Another example might be to predict whether an Indian voter will vote BJP or Trinamool Congress or Left Front or Congress, based on age, income, sex, race, state of residence, votes in previous elections, etc.^[10] The technique can also be used in engineering, especially for predicting the probability of failure of a given process, system or product.^{[11][12]} It is also used in marketing applications such as prediction of a customer's propensity to purchase a product or halt a subscription, etc.^[13] In economics it can be used to predict the likelihood of a person's choosing to be in the labor force, and a business application would be to predict the likelihood of a homeowner defaulting on a mortgage. Conditional random fields, an extension of logistic regression to sequential data, are used in natural language processing.

Examples

Logistic model

Let us try to understand logistic regression by considering a logistic model with given parameters, then seeing how the coefficients can be estimated from data. Consider a model with two predictors, $\boldsymbol{x_1}$ and $\boldsymbol{x_2}$, and one binary (Bernoulli) response variable \boldsymbol{Y} , for which we denote $\boldsymbol{p} = \boldsymbol{P(Y = 1)} = \boldsymbol{E(Y)}$. Predictors may be continuous variables (taking a real number as value), or binary variables (taking value 0 or 1). Then the general form of the log-odds (here denoted by ℓ) is:

$$\ell = \log \frac{p}{1-p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

The log odds are the logarithm of the odds (i.e. the ratio between a probability value and its complementary).

The coefficients β_i are the parameters of the model. Note that this is a linear model: the log-odds ℓ are a linear combination of the predictors $\boldsymbol{x_1}$ and $\boldsymbol{x_2}$, including a constant term β_0 . The corresponding odds are the exponential:

$$\boldsymbol{o} = \boldsymbol{b}^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}$$

where \boldsymbol{b} is the base of the logarithm. This is now a non-linear model since the odds are not a linear combination of the predictors.

By simple algebraic manipulation, the corresponding probability is:

$$\boldsymbol{p} = \frac{\boldsymbol{b}^{\beta_0 + \beta_1 x_1 + \beta_2 x_2}}{\boldsymbol{b}^{\beta_0 + \beta_1 x_1 + \beta_2 x_2} + 1} = \frac{1}{1 + \boldsymbol{b}^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

The base \boldsymbol{b} is usually taken to be \boldsymbol{e} , but for exposition we here use 10, so the odds are more understandable. Suppose the coefficients are $\text{---}\mathbf{3}, \mathbf{1}, \mathbf{2}$, so the model is:

$$\ell = \text{---}\mathbf{3} + \mathbf{1} \cdot \boldsymbol{x_1} + \mathbf{2} \cdot \boldsymbol{x_2},$$

This can be interpreted as follows:

- $\beta_0 = \text{---}\mathbf{3}$ is the y-intercept, and can thus be interpreted as the log-odds when the predictors are all zero; here the log-odds are $\text{---}\mathbf{3}$, so the odds are $\mathbf{10^{-3}} = \mathbf{10^{-3} : 1} = \mathbf{1 : 1000}$, and the probability is $\mathbf{1/(1000 + 1)} = \mathbf{1/1001} \approx \mathbf{0.001} = \mathbf{0.1\%}$.
 - Alternatively, grouping the weighted predictors as a single contribution to the total logit, $\ell' = \mathbf{1} \cdot \boldsymbol{x_1} + \mathbf{2} \cdot \boldsymbol{x_2}$, so $\ell = \ell' - \mathbf{3}$, the quantity $\text{---}\beta_0$ is the x-intercept, and can be interpreted as the weighted predictors corresponding to even odds: if $\ell' = \mathbf{1} \cdot \boldsymbol{x_1} + \mathbf{2} \cdot \boldsymbol{x_2} = \mathbf{3}$, then the log-odds are 0, the odds for are $\mathbf{10^0} = \mathbf{10^0 : 1} = \mathbf{1 : 1}$, and the probability is $\mathbf{1/(1 + 1)} = \mathbf{1/2} = \mathbf{50\%}$.
- $\beta_1 = \mathbf{1}$ means that increasing $\boldsymbol{x_1}$ by 1 increases the log-odds by $\mathbf{1} \cdot \mathbf{1}$, so it *multiplies* the odds by $\mathbf{10^1} = \mathbf{10}$; this is sometimes referred to as the "effect" of the predictor $\boldsymbol{x_1}$.
- $\beta_2 = \mathbf{2}$ means that increasing $\boldsymbol{x_2}$ by 1 increases the log-odds by $\mathbf{2} \cdot \mathbf{1}$, so it *multiplies* the odds by $\mathbf{10^2} = \mathbf{100}$. Thus the effect of $\boldsymbol{x_2}$ on the log-odds is twice as great as the effect of $\boldsymbol{x_1}$.

Such a model can be used for various purposes. For example, given an individual datum with values of predictors $\boldsymbol{x_1}, \boldsymbol{x_2}$, one can estimate the log-odds (hence odds and probability) of the outcome by putting in the values in the formula. Alternatively, consider comparing two medical treatments, one of which decreases $\boldsymbol{x_1}$ by 3 and the other that decreases $\boldsymbol{x_2}$ by 2. Assuming this model is valid, the first treatment reduces the log-odds of the outcome by $\mathbf{1} \cdot \mathbf{3} = \mathbf{3}$, so it reduces the odds by a factor of $\mathbf{10^3} = \mathbf{1000}$, but the other treatment reduces the log-odds by $\mathbf{2} \cdot \mathbf{2} = \mathbf{4}$, so it reduces the odds by a factor of $\mathbf{10^4} = \mathbf{10000}$, and thus is more effective, all else equal.

In order to estimate the parameters of such a logistic model and compute how well it fits the data, one must do logistic regression.

Probability of passing an exam versus hours of study

To answer the following question:

A group of 20 students spends between 0 and 6 hours studying for an exam. How does the number of hours spent studying effect the probability of the student passing the exam?

The reason for using logistic regression for this problem is that the values of the dependent variable, pass and fail, while represented by "1" and "0", are not cardinal numbers. If the problem was changed so that pass/fail was replaced with the grade 0–100 (cardinal numbers), then simple regression analysis could be used.

The table shows the number of hours each student spent studying, and whether they passed (1) or failed (0).

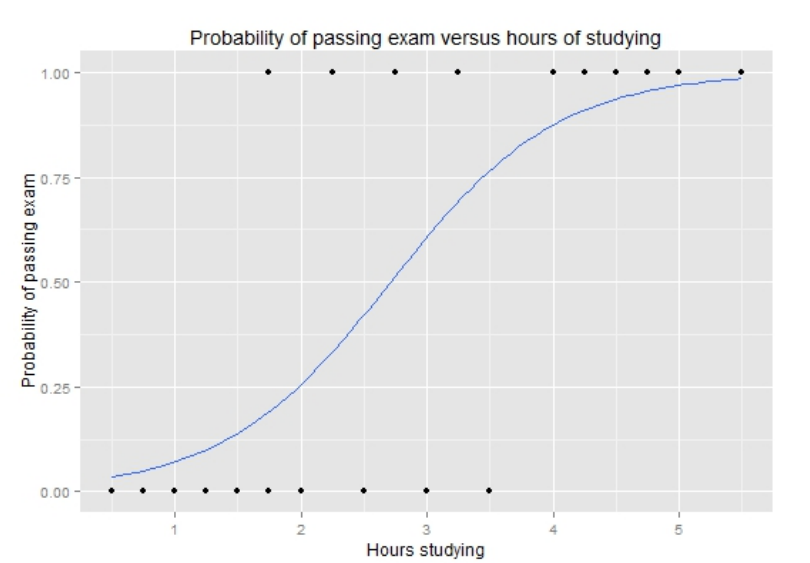
Hours	0.50	0.75	1.00	1.25	1.50	1.75	1.75	2.00	2.25	2.50	2.75	3.00	3.25	3.50	4.00	4.25	4.50	4.75	5.00	5.50
Pass	0	0	0	0	0	0	1	0	1	0	1	0	1	0	1	1	1	1	1	1

The graph shows the probability of passing the exam versus the number of hours studying, with the logistic regression curve fitted to the data.

The logistic regression analysis gives the following output.

	Coefficient	Std.Error	z-value	P-value (Wald)
Intercept	−4.0777	1.7610	−2.316	0.0206
Hours	1.5046	0.6287	2.393	0.0167

The output indicates that hours studying is significantly associated with the probability of passing the exam ($p = \mathbf{0.0167}$, Wald test). The output also provides the coefficients for **Intercept** = −**4.0777** and **Hours** = **1.5046**. These coefficients are entered in the logistic regression equation to estimate the odds (probability) of passing the exam:



Graph of a logistic regression curve showing probability of passing an exam versus hours studying

$$\text{Log-odds of passing exam} = 1.5046 \cdot \text{Hours} - 4.0777 = 1.5046 \cdot (\text{Hours} - 2.71)$$

$$\text{Odds of passing exam} = \exp(1.5046 \cdot \text{Hours} - 4.0777) = \exp(1.5046 \cdot (\text{Hours} - 2.71))$$

$$\text{Probability of passing exam} = \frac{1}{1 + \exp(-(1.5046 \cdot \text{Hours} - 4.0777))}$$

One additional hour of study is estimated to increase log-odds of passing by 1.5046, so multiplying odds of passing by **exp(1.5046) ≈ 4.5**. The form with the x-intercept (2.71) shows that this estimates even odds (log-odds 0, odds 1, probability 1/2) for a student who studies 2.71 hours.

For example, for a student who studies 2 hours, entering the value **Hours = 2** in the equation gives the estimated probability of passing the exam of 0.26:

$$\text{Probability of passing exam} = \frac{1}{1 + \exp(-(1.5046 \cdot 2 - 4.0777))} = \mathbf{0.26}$$

Similarly, for a student who studies 4 hours, the estimated probability of passing the exam is 0.87:

$$\text{Probability of passing exam} = \frac{1}{1 + \exp(-(1.5046 \cdot 4 - 4.0777))} = \mathbf{0.87}$$

This table shows the probability of passing the exam for several values of hours studying.

Hours of study	Passing exam		
	Log-odds	Odds	Probability
1	−2.57	0.076 ≈ 1:13.1	0.07
2	−1.07	0.34 ≈ 1:2.91	0.26
3	0.44	1.55	0.61
4	1.94	6.96	0.87
5	3.45	31.4	0.97

The output from the logistic regression analysis gives a p-value of $p = \mathbf{0.0167}$, which is based on the Wald z-score. Rather than the Wald method, the recommended method to calculate the p-value for logistic regression is the likelihood-ratio test (LRT), which for this data gives $p = \mathbf{0.0006}$.

Discussion

Logistic regression can be binomial, ordinal or multinomial. Binomial or binary logistic regression deals with situations in which the observed outcome for a dependent variable can have only two possible types, "0" and "1" (which may represent, for example, "dead" vs. "alive" or "win" vs. "loss"). Multinomial logistic regression deals with situations where the outcome can have three or more possible types (e.g., "disease A" vs. "disease B" vs. "disease C") that are not ordered. Ordinal logistic regression deals with dependent variables that are ordered.

In binary logistic regression, the outcome is usually coded as "0" or "1", as this leads to the most straightforward interpretation.^[14] If a particular observed outcome for the dependent variable is the noteworthy possible outcome (referred to as a "success" or a "case") it is usually coded as "1" and the contrary outcome (referred to as a "failure" or a "noncase") as "0". Binary logistic regression is used to predict the odds of being a case based on the values of the independent variables (predictors). The odds are defined as the probability that a particular outcome is a case divided by the probability that it is a noncase.

Like other forms of regression analysis, logistic regression makes use of one or more predictor variables that may be either continuous or categorical. Unlike ordinary linear regression, however, logistic regression is used for predicting dependent variables that take membership in one of a limited number of categories (treating the dependent variable in the binomial case as the outcome of a Bernoulli trial) rather than a continuous outcome. Given this difference, the assumptions of linear regression are violated. In particular, the residuals cannot be normally distributed. In addition, linear regression may make nonsensical predictions for a binary dependent variable. What is needed is a way to convert a binary variable into a continuous one that can take on any real value (negative or positive). To do that, binomial logistic regression first calculates the odds of the event happening for different levels of each independent variable, and then takes its logarithm to create a continuous criterion as a transformed version of the dependent variable. The logarithm of the odds is the logit of the probability, the logit is defined as follows:

$$\text{logit } p = \ln \frac{p}{1 - p} \quad \text{for } 0 < p < 1.$$

Although the dependent variable in logistic regression is Bernoulli, the logit is on an unrestricted scale.^[14] The logit function is the link function in this kind of generalized linear model, i.e.

$$\text{logit } E(Y) = \alpha + \beta x$$

Y is the Bernoulli-distributed response variable and x is the predictor variable.

The logit of the probability of success is then fitted to the predictors. The predicted value of the logit is converted back into predicted odds via the inverse of the natural logarithm, namely the exponential function. Thus, although the observed dependent variable in binary logistic regression is a 0-or-1 variable, the logistic regression estimates the odds, as a continuous variable, that the dependent variable is a success (a case). In some applications, the odds are all that is needed. In others, a specific yes-or-no prediction is needed for whether the dependent variable is or is not a case; this categorical prediction can be based on the computed odds of success, with predicted odds above some chosen cutoff value being translated into a prediction of success.

The assumption of linear predictor effects can easily be relaxed using techniques such as spline functions.^[15]

Logistic regression vs. other approaches

Logistic regression measures the relationship between the categorical dependent variable and one or more independent variables by estimating probabilities using a logistic function, which is the cumulative logistic distribution. Thus, it treats the same set of problems as probit regression using similar techniques, with the latter using a cumulative normal distribution curve instead. Equivalently, in the latent variable interpretations of these two methods, the first assumes a standard logistic distribution of errors and the second a standard normal distribution of errors.^[16]

Logistic regression can be seen as a special case of the generalized linear model and thus analogous to linear regression. The model of logistic regression, however, is based on quite different assumptions (about the relationship between the dependent and independent variables) from those of linear regression. In particular, the key differences between these two models can be seen in the following two features of logistic regression. First, the conditional distribution $\boldsymbol{y} \mid \boldsymbol{x}$ is a Bernoulli distribution rather than a Gaussian distribution, because the dependent variable is binary. Second, the predicted values are probabilities and are therefore restricted to (0,1) through the logistic distribution function because logistic regression predicts the **probability** of particular outcomes rather than the outcomes themselves.

Logistic regression is an alternative to Fisher's 1936 method, linear discriminant analysis.^[17] If the assumptions of linear discriminant analysis hold, the conditioning can be reversed to produce logistic regression. The converse is not true, however, because logistic regression does not require the multivariate normal assumption of discriminant analysis.^[18]

Latent variable interpretation

The logistic regression can be understood simply as finding the β parameters that best fit:

$$\boldsymbol{y} = \begin{cases} 1 & \beta_0 + \beta_1 \boldsymbol{x} + \varepsilon > 0 \\ 0 & \text{else} \end{cases}$$

where ε is an error distributed by the standard logistic distribution. (If the standard normal distribution is used instead, it is a probit model.)

The associated latent variable is $\boldsymbol{y}' = \beta_0 + \beta_1 \boldsymbol{x} + \varepsilon$. The error term ε is not observed, and so the \boldsymbol{y}' is also an unobservable, hence termed "latent" (the observed data are values of \boldsymbol{y} and \boldsymbol{x}). Unlike ordinary regression, however, the β parameters cannot be expressed by any direct formula of the \boldsymbol{y} and \boldsymbol{x} values in the observed data. Instead they are to be found by an iterative search process, usually implemented by a software program, that finds the maximum of a complicated "likelihood expression" that is a function of all of the observed \boldsymbol{y} and \boldsymbol{x} values. The estimation approach is explained below.

Logistic function, odds, odds ratio, and logit

Definition of the logistic function

An explanation of logistic regression can begin with an explanation of the standard logistic function. The logistic function is a sigmoid function, which takes any real input t , ($t \in \mathbb{R}$), and outputs a value between zero and one;^[14] for the logit, this is interpreted as taking input log-odds and having output probability. The *standard* logistic function $\sigma : \mathbb{R} \rightarrow (0, 1)$ is defined as follows:

$$\sigma(t) = \frac{e^t}{e^t + 1} = \frac{1}{1 + e^{-t}}$$

A graph of the logistic function on the t -interval $(-6,6)$ is shown in Figure 1.

Let us assume that t is a linear function of a single explanatory variable \boldsymbol{x} (the case where t is a *linear combination* of multiple explanatory variables is treated similarly). We can then express t as follows:

$$t = \beta_0 + \beta_1 \boldsymbol{x}$$

And the general logistic function $\boldsymbol{p} : \mathbb{R} \rightarrow (0, 1)$ can now be written as:

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

In the logistic model, $\boldsymbol{p}(\boldsymbol{x})$ is interpreted as the probability of the dependent variable \boldsymbol{Y} equaling a success/case rather than a failure/non-case. It's clear that the response variables \boldsymbol{Y}_i are not identically distributed: $P(\boldsymbol{Y}_i = 1 \mid \boldsymbol{X})$ differs from one data point \boldsymbol{X}_i to another, though they are independent given design matrix \boldsymbol{X} and shared parameters β .^[8]

Definition of the inverse of the logistic function

We can now define the logit (log odds) function as the inverse $\boldsymbol{g} = \sigma^{-1}$ of the standard logistic function. It is easy to see that it satisfies:

$$\boldsymbol{g}(\boldsymbol{p}(\boldsymbol{x})) = \sigma^{-1}(\boldsymbol{p}(\boldsymbol{x})) = \text{logit } \boldsymbol{p}(\boldsymbol{x}) = \ln\left(\frac{\boldsymbol{p}(\boldsymbol{x})}{1 - \boldsymbol{p}(\boldsymbol{x})}\right) = \beta_0 + \beta_1 \boldsymbol{x},$$

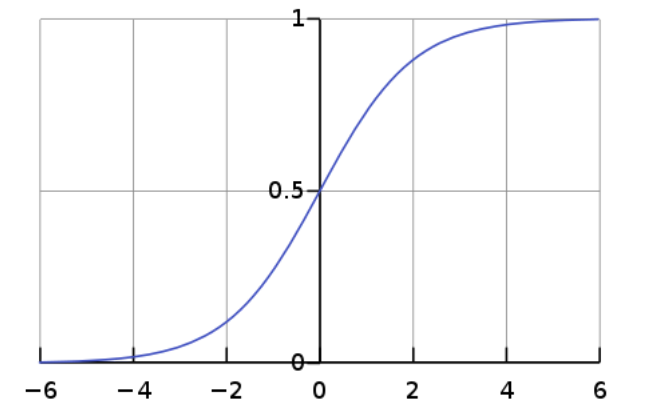


Figure 1. The standard logistic function $\sigma(t)$; note that $\sigma(t) \in (0, 1)$ for all t .

and equivalently, after exponentiating both sides we have the odds:

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

Interpretation of these terms

In the above equations, the terms are as follows:

- g is the logit function. The equation for $g(p(x))$ illustrates that the logit (i.e., log-odds or natural logarithm of the odds) is equivalent to the linear regression expression.
- \ln denotes the natural logarithm.
- $p(x)$ is the probability that the dependent variable equals a case, given some linear combination of the predictors. The formula for $p(x)$ illustrates that the probability of the dependent variable equaling a case is equal to the value of the logistic function of the linear regression expression. This is important in that it shows that the value of the linear regression expression can vary from negative to positive infinity and yet, after transformation, the resulting expression for the probability $p(x)$ ranges between 0 and 1.
- β_0 is the intercept from the linear regression equation (the value of the criterion when the predictor is equal to zero).
- $\beta_1 x$ is the regression coefficient multiplied by some value of the predictor.
- base e denotes the exponential function.

Definition of the odds

The odds of the dependent variable equaling a case (given some linear combination x of the predictors) is equivalent to the exponential function of the linear regression expression. This illustrates how the logit serves as a link function between the probability and the linear regression expression. Given that the logit ranges between negative and positive infinity, it provides an adequate criterion upon which to conduct linear regression and the logit is easily converted back into the odds.^[14]

So we define odds of the dependent variable equaling a case (given some linear combination x of the predictors) as follows:

$$\text{odds} = e^{\beta_0 + \beta_1 x}.$$

The odds ratio

For a continuous independent variable the odds ratio can be defined as:

$$\text{OR} = \frac{\text{odds}(x + 1)}{\text{odds}(x)} = \frac{\left(\frac{F(x+1)}{1-F(x+1)}\right)}{\left(\frac{F(x)}{1-F(x)}\right)} = \frac{e^{\beta_0 + \beta_1 (x+1)}}{e^{\beta_0 + \beta_1 x}} = e^{\beta_1}$$

This exponential relationship provides an interpretation for β_1 : The odds multiply by e^{β_1} for every 1-unit increase in x.^[19]

For a binary independent variable the odds ratio is defined as $\frac{ad}{bc}$ where a, b, c and d are cells in a 2×2 contingency table.^[20]

Multiple explanatory variables

If there are multiple explanatory variables, the above expression $\beta_0 + \beta_1 x$ can be revised to $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_m x_m = \beta_0 + \sum_{i=1}^m \beta_i x_i$. Then when this is used in the equation relating the log odds of a success to the values of the predictors, the linear regression will be a multiple regression with m explanators; the parameters β_j for all $j = 0, 1, 2, \dots, m$ are all estimated.

Again, the more traditional equations are:

$$\log \frac{p}{1 - p} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_m x_m$$

and

$$p = \frac{1}{1 + b^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_m x_m)}}$$

where usually $b = e$.

Model fitting

Logistic regression is an important machine learning algorithm. The goal is to model the probability of a random variable Y being 0 or 1 given experimental data.^[21]

Consider a generalized linear model function parameterized by θ ,

$$h_{\theta}(X) = \frac{1}{1 + e^{-\theta^T X}} = \text{Pr}(Y = 1 \mid X; \theta)$$

Therefore,

$$\text{Pr}(Y = 0 \mid X; \theta) = 1 - h_{\theta}(X)$$

and since $Y \in \{0, 1\}$, $\text{Pr}(y \mid X; \theta)$ is given by the following equation using an exponentiation trick.

If we attempt to model the probability that \boldsymbol{Y} is 0 or 1 with the function $\Pr(\boldsymbol{y} \mid \boldsymbol{X}; \boldsymbol{\theta}) = h_{\boldsymbol{\theta}}(\boldsymbol{X})^{\boldsymbol{y}}(1 - h_{\boldsymbol{\theta}}(\boldsymbol{X}))^{(1-\boldsymbol{y})}$, we take our likelihood function assuming that all the observations in the sample are independently Bernoulli distributed,

$$\begin{aligned} L(\boldsymbol{\theta} \mid \boldsymbol{x}) &= \Pr(\boldsymbol{Y} \mid \boldsymbol{X}; \boldsymbol{\theta}) \\ &= \prod_i \Pr(y_i \mid x_i; \boldsymbol{\theta}) \\ &= \prod_i h_{\boldsymbol{\theta}}(x_i)^{y_i} (1 - h_{\boldsymbol{\theta}}(x_i))^{(1-y_i)} \end{aligned}$$

Typically, the log likelihood is maximized with a normalizing factor N^{-1} ,

$$N^{-1} \log L(\boldsymbol{\theta} \mid \boldsymbol{x}) = N^{-1} \sum_{i=1}^N \log \Pr(y_i \mid x_i; \boldsymbol{\theta})$$

which is maximized using optimization techniques such as gradient descent.

Assuming the $(\boldsymbol{x}, \boldsymbol{y})$ pairs are drawn uniformly from the underlying distribution, then in the limit of large N ,

$$\begin{aligned} \lim_{N \rightarrow +\infty} N^{-1} \sum_{i=1}^N \log \Pr(y_i \mid x_i; \boldsymbol{\theta}) &= \sum_{\boldsymbol{x} \in \mathcal{X}} \sum_{\boldsymbol{y} \in \mathcal{Y}} \Pr(\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Y} = \boldsymbol{y}) \log \Pr(\boldsymbol{Y} = \boldsymbol{y} \mid \boldsymbol{X} = \boldsymbol{x}; \boldsymbol{\theta}) \\ &= \sum_{\boldsymbol{x} \in \mathcal{X}} \sum_{\boldsymbol{y} \in \mathcal{Y}} \Pr(\boldsymbol{X} = \boldsymbol{x}, \boldsymbol{Y} = \boldsymbol{y}) \left(-\log \frac{\Pr(\boldsymbol{Y} = \boldsymbol{y} \mid \boldsymbol{X} = \boldsymbol{x})}{\Pr(\boldsymbol{Y} = \boldsymbol{y} \mid \boldsymbol{X} = \boldsymbol{x}; \boldsymbol{\theta})} + \log \Pr(\boldsymbol{Y} = \boldsymbol{y} \mid \boldsymbol{X} = \boldsymbol{x}) \right) \\ &= -D_{\text{KL}}(\boldsymbol{Y} \parallel \boldsymbol{Y}_{\boldsymbol{\theta}}) - H(\boldsymbol{Y} \mid \boldsymbol{X}) \end{aligned}$$

where $H(\boldsymbol{X} \mid \boldsymbol{Y})$ is the conditional entropy and D_{KL} is the Kullback–Leibler divergence. This leads to the intuition that by maximizing the log-likelihood of a model, you are minimizing the KL divergence of your model from the maximal entropy distribution. Intuitively searching for the model that makes the fewest assumptions in its parameters.

"Rule of ten"

A widely used rule of thumb, the "one in ten rule", states that logistic regression models give stable values for the explanatory variables if based on a minimum of about 10 events per explanatory variable (EPV); where *event* denotes the cases belonging to the less frequent category in the dependent variable. Thus a study designed to use \boldsymbol{k} explanatory variables for an event (e.g. myocardial infarction) expected to occur in a proportion \boldsymbol{p} of participants in the study will require a total of $\mathbf{10k/p}$ participants. However, there is considerable debate about the reliability of this rule, which is based on simulation studies and lacks a secure theoretical underpinning.^[22] According to some authors^[23] the rule is overly conservative, some circumstances; with the authors stating "If we (somewhat subjectively) regard confidence interval coverage less than 93 percent, type I error greater than 7 percent, or relative bias greater than 15 percent as problematic, our results indicate that problems are fairly frequent with 2–4 EPV, uncommon with 5–9 EPV, and still observed with 10–16 EPV. The worst instances of each problem were not severe with 5–9 EPV and usually comparable to those with 10–16 EPV".^[24]

Others have found results that are not consistent with the above, using different criteria. A useful criterion is whether the fitted model will be expected to achieve the same predictive discrimination in a new sample as it appeared to achieve in the model development sample. For that criterion, 20 events per candidate variable may be required.^[25] Also, one can argue that 96 observations are needed only to estimate the model's intercept precisely enough that the margin of error in predicted probabilities is ± 0.1 with an 0.95 confidence level.^[15]

Maximum likelihood estimation

The regression coefficients are usually estimated using maximum likelihood estimation.^[26] Unlike linear regression with normally distributed residuals, it is not possible to find a closed-form expression for the coefficient values that maximize the likelihood function, so that an iterative process must be used instead; for example Newton's method. This process begins with a tentative solution, revises it slightly to see if it can be improved, and repeats this revision until no more improvement is made, at which point the process is said to have converged.^[26]

In some instances, the model may not reach convergence. Non-convergence of a model indicates that the coefficients are not meaningful because the iterative process was unable to find appropriate solutions. A failure to converge may occur for a number of reasons: having a large ratio of predictors to cases, multicollinearity, sparseness, or complete separation.

- Having a large ratio of variables to cases results in an overly conservative Wald statistic (discussed below) and can lead to non-convergence.
- Multicollinearity refers to unacceptably high correlations between predictors. As multicollinearity increases, coefficients remain unbiased but standard errors increase and the likelihood of model convergence decreases.^[26] To detect multicollinearity amongst the predictors, one can conduct a linear regression analysis with the predictors of interest for the sole purpose of examining the tolerance statistic ^[26] used to assess whether multicollinearity is unacceptably high.
- Sparseness in the data refers to having a large proportion of empty cells (cells with zero counts). Zero cell counts are particularly problematic with categorical predictors. With continuous predictors, the model can infer values for the zero cell counts, but this is not the case with categorical predictors. The model will not converge with zero cell counts for categorical predictors because the natural logarithm of zero is an undefined value so that the final solution to the model cannot be reached. To remedy this problem, researchers may collapse categories in a theoretically meaningful way or add a constant to all cells.^[26]
- Another numerical problem that may lead to a lack of convergence is complete separation, which refers to the instance in which the predictors perfectly predict the criterion – all cases are accurately classified. In such instances, one should reexamine the data, as there is likely some kind of error.^[14]

Iteratively reweighted least squares (IRLS)

Binary logistic regression ($\boldsymbol{y} = \mathbf{0}$ or $\boldsymbol{y} = \mathbf{1}$) can, for example, be calculated using *iteratively reweighted least squares* (IRLS), which is equivalent to minimizing the log-likelihood of a Bernoulli distributed process using Newton's method. If the problem is written in vector matrix form, with parameters $\boldsymbol{w}^T = [\beta_0, \beta_1, \beta_2, \dots]$, explanatory variables $\boldsymbol{x}(i) = [\mathbf{1}, x_1(i), x_2(i), \dots]^T$ and expected value of the Bernoulli distribution $\boldsymbol{\mu}(i) = \frac{1}{1 + e^{-\boldsymbol{w}^T \boldsymbol{x}(i)}}$, the parameters \boldsymbol{w} can be found using the following iterative algorithm:

$\mathbf{w}_{k+1} = \left(\mathbf{X}^T \mathbf{S}_k \mathbf{X}\right)^{-1} \mathbf{X}^T \left(\mathbf{S}_k \mathbf{X} \mathbf{w}_k + \mathbf{y} - \boldsymbol{\mu}_k\right)$

where $\mathbf{S} = \text{diag}(\mu(i)(1 - \mu(i)))$ is a diagonal weighting matrix, $\boldsymbol{\mu} = [\mu(1), \mu(2), \dots]$ the vector of expected values,

$$\mathbf{X} = \begin{bmatrix} 1 & x_1(1) & x_2(1) & \dots \\ 1 & x_1(2) & x_2(2) & \dots \\ \vdots & \vdots & \vdots & \end{bmatrix}$$

The regressor matrix and $\mathbf{y}(i) = [y(1), y(2), \dots]^T$ the vector of response variables. More details can be found e.g. here ^[27]

Evaluating goodness of fit

Goodness of fit in linear regression models is generally measured using R². Since this has no direct analog in logistic regression, various methods^{[28]:ch.21} including the following can be used instead.

Deviance and likelihood ratio tests

In linear regression analysis, one is concerned with partitioning variance via the sum of squares calculations – variance in the criterion is essentially divided into variance accounted for by the predictors and residual variance. In logistic regression analysis, deviance is used in lieu of a sum of squares calculations.^[29] Deviance is analogous to the sum of squares calculations in linear regression^[14] and is a measure of the lack of fit to the data in a logistic regression model.^[29] When a "saturated" model is available (a model with a theoretically perfect fit), deviance is calculated by comparing a given model with the saturated model.^[14] This computation gives the likelihood-ratio test.^[14]

$$D = -2 \ln \frac{\text{likelihood of the fitted model}}{\text{likelihood of the saturated model}}.$$

In the above equation, *D* represents the deviance and ln represents the natural logarithm. The log of this likelihood ratio (the ratio of the fitted model to the saturated model) will produce a negative value, hence the need for a negative sign. *D* can be shown to follow an approximate chi-squared distribution.^[14] Smaller values indicate better fit as the fitted model deviates less from the saturated model. When assessed upon a chi-square distribution, nonsignificant chi-square values indicate very little unexplained variance and thus, good model fit. Conversely, a significant chi-square value indicates that a significant amount of the variance is unexplained.

When the saturated model is not available (a common case), deviance is calculated simply as −2·(log likelihood of the fitted model), and the reference to the saturated model's log likelihood can be removed from all that follows without harm.

Two measures of deviance are particularly important in logistic regression: null deviance and model deviance. The null deviance represents the difference between a model with only the intercept (which means "no predictors") and the saturated model. The model deviance represents the difference between a model with at least one predictor and the saturated model.^[29] In this respect, the null model provides a baseline upon which to compare predictor models. Given that deviance is a measure of the difference between a given model and the saturated model, smaller values indicate better fit. Thus, to assess the contribution of a predictor or set of predictors, one can subtract the model deviance from the null deviance and assess the difference on a χ^2_{s-p} , chi-square distribution with degrees of freedom^[14] equal to the difference in the number of parameters estimated.

Let

$$D_{\text{null}} = -2 \ln \frac{\text{likelihood of null model}}{\text{likelihood of the saturated model}}$$

$$D_{\text{fitted}} = -2 \ln \frac{\text{likelihood of fitted model}}{\text{likelihood of the saturated model}}.$$

Then the difference of both is:

$$\begin{aligned} D_{\text{null}} - D_{\text{fitted}} &= -2 \left(\ln \frac{\text{likelihood of null model}}{\text{likelihood of the saturated model}} - \ln \frac{\text{likelihood of fitted model}}{\text{likelihood of the saturated model}} \right) \\ &= -2 \ln \frac{\left(\frac{\text{likelihood of null model}}{\text{likelihood of the saturated model}} \right)}{\left(\frac{\text{likelihood of fitted model}}{\text{likelihood of the saturated model}} \right)} \\ &= -2 \ln \frac{\text{likelihood of the null model}}{\text{likelihood of fitted model}}. \end{aligned}$$

If the model deviance is significantly smaller than the null deviance then one can conclude that the predictor or set of predictors significantly improved model fit. This is analogous to the *F*-test used in linear regression analysis to assess the significance of prediction.^[29]

Pseudo-*R*²s

In linear regression the squared multiple correlation, *R*² is used to assess goodness of fit as it represents the proportion of variance in the criterion that is explained by the predictors.^[29] In logistic regression analysis, there is no agreed upon analogous measure, but there are several competing measures each with limitations.^{[29][30]}

Four of the most commonly used indices and one less commonly used one are examined on this page:

- Likelihood ratio *R*²_L
- Cox and Snell *R*²_{CS}
- Nagelkerke *R*²_N

- McFadden R^2_{McF}
- Tjur R^2_{T}

R^2_{L} is given by ^[29]

$$R^2_{\text{L}} = \frac{D_{\text{null}} - D_{\text{fitted}}}{D_{\text{null}}}.$$

This is the most analogous index to the squared multiple correlations in linear regression.^[26] It represents the proportional reduction in the deviance wherein the deviance is treated as a measure of variation analogous but not identical to the variance in linear regression analysis.^[26] One limitation of the likelihood ratio R^2 is that it is not monotonically related to the odds ratio,^[29] meaning that it does not necessarily increase as the odds ratio increases and does not necessarily decrease as the odds ratio decreases.

R^2_{CS} is an alternative index of goodness of fit related to the R^2 value from linear regression.^[30] It is given by:

$$\begin{aligned} R^2_{\text{CS}} &= 1 - \left(\frac{L_0}{L_M}\right)^{2/n} \\ &= 1 - e^{2(\ln(L_0) - \ln(L_M))/n} \end{aligned}$$

where L_M and L_0 are the likelihoods for the model being fitted and the null model, respectively. The Cox and Snell index is problematic as its maximum value is $1 - L_0^{2/n}$. The highest this upper bound can be is 0.75, but it can easily be as low as 0.48 when the marginal proportion of cases is small.^[30]

R^2_{N} provides a correction to the Cox and Snell R^2 so that the maximum value is equal to 1. Nevertheless, the Cox and Snell and likelihood ratio R^2 s show greater agreement with each other than either does with the Nagelkerke R^2 .^[29] Of course, this might not be the case for values exceeding .75 as the Cox and Snell index is capped at this value. The likelihood ratio R^2 is often preferred to the alternatives as it is most analogous to R^2 in linear regression, is independent of the base rate (both Cox and Snell and Nagelkerke R^2 s increase as the proportion of cases increase from 0 to .5) and varies between 0 and 1.

R^2_{McF} is defined as

$$R^2_{\text{McF}} = 1 - \frac{\ln(L_0)}{\ln L_M},$$

and is preferred over R^2_{CS} by Allison.^[30] The two expressions R^2_{McF} and R^2_{CS} are then related respectively by,

$$R^2_{\text{CS}} = 1 - \left(\frac{1}{L_0}\right)^{\frac{2(R^2_{\text{McF}})}{n}}$$

$$R^2_{\text{McF}} = -\frac{n}{2} \cdot \frac{\ln(1 - R^2_{\text{CS}})}{\ln L_0}$$

However, Allison now prefers R^2_{T} which is a relatively new measure developed by Tjur.^[31] It can be calculated in two steps:^[30]

- For each level of the dependent variable, find the mean of the predicted probabilities of an event.
- Take the absolute value of the difference between these means

A word of caution is in order when interpreting pseudo- R^2 statistics. The reason these indices of fit are referred to as *pseudo* R^2 is that they do not represent the proportionate reduction in error as the R^2 in linear regression does.^[29] Linear regression assumes homoscedasticity, that the error variance is the same for all values of the criterion. Logistic regression will always be heteroscedastic – the error variances differ for each value of the predicted score. For each value of the predicted score there would be a different value of the proportionate reduction in error. Therefore, it is inappropriate to think of R^2 as a proportionate reduction in error in a universal sense in logistic regression.^[29]

Hosmer–Lemeshow test

The Hosmer–Lemeshow test uses a test statistic that asymptotically follows a χ^2 distribution to assess whether or not the observed event rates match expected event rates in subgroups of the model population. This test is considered to be obsolete by some statisticians because of its dependence on arbitrary binning of predicted probabilities and relative low power.^[32]

Coefficients

After fitting the model, it is likely that researchers will want to examine the contribution of individual predictors. To do so, they will want to examine the regression coefficients. In linear regression, the regression coefficients represent the change in the criterion for each unit change in the predictor.^[29] In logistic regression, however, the regression coefficients represent the change in the logit for each unit change in the predictor. Given that the logit is not intuitive, researchers are likely to focus on a predictor's effect on the exponential function of the regression coefficient – the odds ratio (see definition). In linear regression, the significance of a regression coefficient is assessed by computing a *t* test. In logistic regression, there are several different tests designed to assess the significance of an individual predictor, most notably the likelihood ratio test and the Wald statistic.

Likelihood ratio test

The likelihood-ratio test discussed above to assess model fit is also the recommended procedure to assess the contribution of individual "predictors" to a given model.^{[14][26][29]} In the case of a single predictor model, one simply compares the deviance of the predictor model with that of the null model on a chi-square distribution with a single degree of freedom. If the predictor model has significantly smaller deviance (c.f chi-square using the difference in degrees of freedom of the two models), then one can conclude that there is a significant association between the "predictor" and the outcome. Although some common statistical packages (e.g. SPSS) do provide likelihood ratio test statistics, without this computationally intensive test it would be more difficult to assess the contribution of individual predictors in the multiple logistic regression case. To assess the contribution of individual predictors one can enter the predictors hierarchically, comparing each new model with the previous to determine

the contribution of each predictor.^[29] There is some debate among statisticians about the appropriateness of so-called "stepwise" procedures. The fear is that they may not preserve nominal statistical properties and may become misleading.^[1] (https://www.amazon.com/Regression-Modeling-Strategies-Applications-Statistics/dp/1441929185/ref=sr_1_2?ie=UTF8&qid=1339171287&sr=8-2)

Wald statistic

Alternatively, when assessing the contribution of individual predictors in a given model, one may examine the significance of the Wald statistic. The Wald statistic, analogous to the *t*-test in linear regression, is used to assess the significance of coefficients. The Wald statistic is the ratio of the square of the regression coefficient to the square of the standard error of the coefficient and is asymptotically distributed as a chi-square distribution.^[26]

W_j = \frac{\beta_j^2}{SE_{\beta_j}^2}

Although several statistical packages (e.g., SPSS, SAS) report the Wald statistic to assess the contribution of individual predictors, the Wald statistic has limitations. When the regression coefficient is large, the standard error of the regression coefficient also tends to be larger increasing the probability of Type-II error. The Wald statistic also tends to be biased when data are sparse.^[29]

Case-control sampling

Suppose cases are rare. Then we might wish to sample them more frequently than their prevalence in the population. For example, suppose there is a disease that affects 1 person in 10,000 and to collect our data we need to do a complete physical. It may be too expensive to do thousands of physicals of healthy people in order to obtain data for only a few diseased individuals. Thus, we may evaluate more diseased individuals, perhaps all of the rare outcomes. This is also retrospective sampling, or equivalently it is called unbalanced data. As a rule of thumb, sampling controls at a rate of five times the number of cases will produce sufficient control data.^[33]

Logistic regression is unique in that it may be estimated on unbalanced data, rather than randomly sampled data, and still yield correct coefficient estimates of the effects of each independent variable on the outcome. That is to say, if we form a logistic model from such data, if the model is correct in the general population, the *β_j* parameters are all correct except for *β₀*. We can correct *β₀* if we know the true prevalence as follows:^[33]

\hat{\beta}_0^* = \hat{\beta}_0 + \log \frac{\pi}{1 - \pi} - \log \frac{\tilde{\pi}}{1 - \tilde{\pi}}

where *π* is the true prevalence and *π̃* is the prevalence in the sample.

Formal mathematical specification

There are various equivalent specifications of logistic regression, which fit into different types of more general models. These different specifications allow for different sorts of useful generalizations.

Setup

The basic setup of logistic regression is as follows. We are given a dataset containing *N* points. Each point *i* consists of a set of *m* input variables *x_{1,i}* ... *x_{m,i}* (also called independent variables, predictor variables, features, or attributes), and a binary outcome variable *Y_i* (also known as a dependent variable, response variable, output variable, or class), i.e. it can assume only the two possible values 0 (often meaning "no" or "failure") or 1 (often meaning "yes" or "success"). The goal of logistic regression is to use the dataset to create a predictive model of the outcome variable.

Some examples:

- The observed outcomes are the presence or absence of a given disease (e.g. diabetes) in a set of patients, and the explanatory variables might be characteristics of the patients thought to be pertinent (sex, race, age, blood pressure, body-mass index, etc.).
- The observed outcomes are the votes (e.g. Democratic or Republican) of a set of people in an election, and the explanatory variables are the demographic characteristics of each person (e.g. sex, race, age, income, etc.). In such a case, one of the two outcomes is arbitrarily coded as 1, and the other as 0.

As in linear regression, the outcome variables *Y_i* are assumed to depend on the explanatory variables *x_{1,i}* ... *x_{m,i}*.

Explanatory variables

As shown above in the above examples, the explanatory variables may be of any type: real-valued, binary, categorical, etc. The main distinction is between continuous variables (such as income, age and blood pressure) and discrete variables (such as sex or race). Discrete variables referring to more than two possible choices are typically coded using dummy variables (or indicator variables), that is, separate explanatory variables taking the value 0 or 1 are created for each possible value of the discrete variable, with a 1 meaning "variable does have the given value" and a 0 meaning "variable does not have that value". For example, a four-way discrete variable of blood type with the possible values "A, B, AB, O" can be converted to four separate two-way dummy variables, "is-A, is-B, is-AB, is-O", where only one of them has the value 1 and all the rest have the value 0. This allows for separate regression coefficients to be matched for each possible value of the discrete variable. (In a case like this, only three of the four dummy variables are independent of each other, in the sense that once the values of three of the variables are known, the fourth is automatically determined. Thus, it is necessary to encode only three of the four possibilities as dummy variables. This also means that when all four possibilities are encoded, the overall model is not identifiable in the absence of additional constraints such as a regularization constraint. Theoretically, this could cause problems, but in reality almost all logistic regression models are fitted with regularization constraints.)

Outcome variables

Formally, the outcomes *Y_i* are described as being Bernoulli-distributed data, where each outcome is determined by an unobserved probability *p_i* that is specific to the outcome at hand, but related to the explanatory variables. This can be expressed in any of the following equivalent forms:

$$Y_i \mid x_{1,i}, \dots, x_{m,i} \sim \text{Bernoulli}(p_i)$$

$$\mathbf{E}[Y_i \mid x_{1,i}, \dots, x_{m,i}] = p_i$$

$$\Pr(Y_i = y \mid x_{1,i}, \dots, x_{m,i}) = \begin{cases} p_i & \text{if } y = 1 \\ 1 - p_i & \text{if } y = 0 \end{cases}$$

$$\Pr(Y_i = y \mid x_{1,i}, \dots, x_{m,i}) = p_i^y (1 - p_i)^{(1-y)}$$

The meanings of these four lines are:

- The first line expresses the probability distribution of each Y_i : Conditioned on the explanatory variables, it follows a Bernoulli distribution with parameters p_i , the probability of the outcome of 1 for trial i . As noted above, each separate trial has its own probability of success, just as each trial has its own explanatory variables. The probability of success p_i is not observed, only the outcome of an individual Bernoulli trial using that probability.
- The second line expresses the fact that the expected value of each Y_i is equal to the probability of success p_i , which is a general property of the Bernoulli distribution. In other words, if we run a large number of Bernoulli trials using the same probability of success p_i , then take the average of all the 1 and 0 outcomes, then the result would be close to p_i . This is because doing an average this way simply computes the proportion of successes seen, which we expect to converge to the underlying probability of success.
- The third line writes out the probability mass function of the Bernoulli distribution, specifying the probability of seeing each of the two possible outcomes.
- The fourth line is another way of writing the probability mass function, which avoids having to write separate cases and is more convenient for certain types of calculations. This relies on the fact that Y_i can take only the value 0 or 1. In each case, one of the exponents will be 1, "choosing" the value under it, while the other is 0, "canceling out" the value under it. Hence, the outcome is either p_i or $1 - p_i$, as in the previous line.

Linear predictor function

The basic idea of logistic regression is to use the mechanism already developed for linear regression by modeling the probability p_i using a linear predictor function, i.e. a linear combination of the explanatory variables and a set of regression coefficients that are specific to the model at hand but the same for all trials. The linear predictor function $f(i)$ for a particular data point i is written as:

$$f(i) = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_m x_{m,i},$$

where β_0, \dots, β_m are regression coefficients indicating the relative effect of a particular explanatory variable on the outcome.

The model is usually put into a more compact form as follows:

- The regression coefficients $\beta_0, \beta_1, ..., \beta_m$ are grouped into a single vector $\boldsymbol{\beta}$ of size $m + 1$.
- For each data point i , an additional explanatory pseudo-variable $x_{0,i}$ is added, with a fixed value of 1, corresponding to the intercept coefficient β_0 .
- The resulting explanatory variables $x_{0,i}, x_{1,i}, ..., x_{m,i}$ are then grouped into a single vector \mathbf{X}_i of size $m + 1$.

This makes it possible to write the linear predictor function as follows:

$$f(i) = \boldsymbol{\beta} \cdot \mathbf{X}_i,$$

using the notation for a dot product between two vectors.

As a generalized linear model

The particular model used by logistic regression, which distinguishes it from standard linear regression and from other types of regression analysis used for binary-valued outcomes, is the way the probability of a particular outcome is linked to the linear predictor function:

$$\text{logit}(\mathbf{E}[Y_i \mid x_{1,i}, \dots, x_{m,i}]) = \text{logit}(p_i) = \ln\left(\frac{p_i}{1 - p_i}\right) = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_m x_{m,i}$$

Written using the more compact notation described above, this is:

$$\text{logit}(\mathbf{E}[Y_i \mid \mathbf{X}_i]) = \text{logit}(p_i) = \ln\left(\frac{p_i}{1 - p_i}\right) = \boldsymbol{\beta} \cdot \mathbf{X}_i$$

This formulation expresses logistic regression as a type of generalized linear model, which predicts variables with various types of probability distributions by fitting a linear predictor function of the above form to some sort of arbitrary transformation of the expected value of the variable.

The intuition for transforming using the logit function (the natural log of the odds) was explained above. It also has the practical effect of converting the probability (which is bounded to be between 0 and 1) to a variable that ranges over $(-\infty, +\infty)$ — thereby matching the potential range of the linear prediction function on the right side of the equation.

Note that both the probabilities p_i and the regression coefficients are unobserved, and the means of determining them is not part of the model itself. They are typically determined by some sort of optimization procedure, e.g. maximum likelihood estimation, that finds values that best fit the observed data (i.e. that give the most accurate predictions for the data already observed), usually subject to regularization conditions that seek to exclude unlikely values, e.g. extremely large values for any of the regression coefficients. The use of a regularization condition is equivalent to doing maximum a posteriori (MAP) estimation, an extension of maximum likelihood. (Regularization is most commonly done using a squared regularizing function, which is equivalent to placing a zero-mean Gaussian prior distribution on the coefficients, but other regularizers are also possible.) Whether or not regularization is used, it is usually not possible to find a closed-form solution; instead, an iterative numerical method must be used, such as iteratively reweighted least squares (IRLS) or, more commonly these days, a quasi-Newton method such as the L-BFGS method.^[34]

The interpretation of the β_j parameter estimates is as the additive effect on the log of the odds for a unit change in the j the explanatory variable. In the case of a dichotomous explanatory variable, for instance, gender e^β is the estimate of the odds of having the outcome for, say, males compared with females.

An equivalent formula uses the inverse of the logit function, which is the logistic function, i.e.:

$$\mathbf{E}[Y_i \mid \mathbf{X}_i] = p_i = \text{logit}^{-1}(\boldsymbol{\beta} \cdot \mathbf{X}_i) = \frac{1}{1 + e^{-\boldsymbol{\beta} \cdot \mathbf{X}_i}}$$

The formula can also be written as a probability distribution (specifically, using a probability mass function):

$$\Pr(Y_i = y \mid \mathbf{X}_i) = p_i^y (1 - p_i)^{1-y} = \left(\frac{e^{\boldsymbol{\beta} \cdot \mathbf{X}_i}}{1 + e^{\boldsymbol{\beta} \cdot \mathbf{X}_i}} \right)^y \left(1 - \frac{e^{\boldsymbol{\beta} \cdot \mathbf{X}_i}}{1 + e^{\boldsymbol{\beta} \cdot \mathbf{X}_i}} \right)^{1-y} = \frac{e^{\boldsymbol{\beta} \cdot \mathbf{X}_i \cdot y}}{1 + e^{\boldsymbol{\beta} \cdot \mathbf{X}_i}}$$

As a latent-variable model

The above model has an equivalent formulation as a latent-variable model. This formulation is common in the theory of discrete choice models and makes it easier to extend to certain more complicated models with multiple, correlated choices, as well as to compare logistic regression to the closely related probit model.

Imagine that, for each trial i , there is a continuous latent variable Y_i^* (i.e. an unobserved random variable) that is distributed as follows:

$$Y_i^* = \boldsymbol{\beta} \cdot \mathbf{X}_i + \varepsilon$$

where

$$\varepsilon \sim \text{Logistic}(0, 1)$$

i.e. the latent variable can be written directly in terms of the linear predictor function and an additive random error variable that is distributed according to a standard logistic distribution.

Then Y_i can be viewed as an indicator for whether this latent variable is positive:

$$Y_i = \begin{cases} 1 & \text{if } Y_i^* > 0 \text{ i.e. } -\varepsilon < \boldsymbol{\beta} \cdot \mathbf{X}_i, \\ 0 & \text{otherwise.} \end{cases}$$

The choice of modeling the error variable specifically with a standard logistic distribution, rather than a general logistic distribution with the location and scale set to arbitrary values, seems restrictive, but in fact, it is not. It must be kept in mind that we can choose the regression coefficients ourselves, and very often can use them to offset changes in the parameters of the error variable's distribution. For example, a logistic error-variable distribution with a non-zero location parameter μ (which sets the mean) is equivalent to a distribution with a zero location parameter, where μ has been added to the intercept coefficient. Both situations produce the same value for Y_i^* regardless of settings of explanatory variables. Similarly, an arbitrary scale parameter s is equivalent to setting the scale parameter to 1 and then dividing all regression coefficients by s . In the latter case, the resulting value of Y_i^* will be smaller by a factor of s than in the former case, for all sets of explanatory variables — but critically, it will always remain on the same side of 0, and hence lead to the same Y_i choice.

(Note that this predicts that the irrelevancy of the scale parameter may not carry over into more complex models where more than two choices are available.)

It turns out that this formulation is exactly equivalent to the preceding one, phrased in terms of the generalized linear model and without any latent variables. This can be shown as follows, using the fact that the cumulative distribution function (CDF) of the standard logistic distribution is the logistic function, which is the inverse of the logit function, i.e.

$$\Pr(\varepsilon < x) = \text{logit}^{-1}(x)$$

Then:

$$\begin{aligned} \Pr(Y_i = 1 \mid \mathbf{X}_i) &= \Pr(Y_i^* > 0 \mid \mathbf{X}_i) \\ &= \Pr(\boldsymbol{\beta} \cdot \mathbf{X}_i + \varepsilon > 0) \\ &= \Pr(\varepsilon > -\boldsymbol{\beta} \cdot \mathbf{X}_i) \\ &= \Pr(\varepsilon < \boldsymbol{\beta} \cdot \mathbf{X}_i) && \text{(because the logistic distribution is symmetric)} \\ &= \text{logit}^{-1}(\boldsymbol{\beta} \cdot \mathbf{X}_i) \\ &= p_i && \text{(see above)} \end{aligned}$$

This formulation—which is standard in discrete choice models—makes clear the relationship between logistic regression (the "logit model") and the probit model, which uses an error variable distributed according to a standard normal distribution instead of a standard logistic distribution. Both the logistic and normal distributions are symmetric with a basic unimodal, "bell curve" shape. The only difference is that the logistic distribution has somewhat heavier tails, which means that it is less sensitive to outlying data (and hence somewhat more robust to model mis-specifications or erroneous data).

Two-way latent-variable model

Yet another formulation uses two separate latent variables:

$$\begin{aligned} Y_i^{0*} &= \boldsymbol{\beta}_0 \cdot \mathbf{X}_i + \varepsilon_0 \\ Y_i^{1*} &= \boldsymbol{\beta}_1 \cdot \mathbf{X}_i + \varepsilon_1 \end{aligned}$$

where

$$\begin{aligned} \varepsilon_0 &\sim \text{EV}_1(0, 1) \\ \varepsilon_1 &\sim \text{EV}_1(0, 1) \end{aligned}$$

where $EV_1(0,1)$ is a standard type-1 extreme value distribution: i.e.

$$\Pr(\varepsilon_0 = x) = \Pr(\varepsilon_1 = x) = e^{-x} e^{-e^{-x}}$$

Then

$$Y_i = \begin{cases} 1 & \text{if } Y_i^{1*} > Y_i^{0*}, \\ 0 & \text{otherwise.} \end{cases}$$

This model has a separate latent variable and a separate set of regression coefficients for each possible outcome of the dependent variable. The reason for this separation is that it makes it easy to extend logistic regression to multi-outcome categorical variables, as in the multinomial logit model. In such a model, it is natural to model each possible outcome using a different set of regression coefficients. It is also possible to motivate each of the separate latent variables as the theoretical utility associated with making the associated choice, and thus motivate logistic regression in terms of utility theory. (In terms of utility theory, a rational actor always chooses the choice with the greatest associated utility.) This is the approach taken by economists when formulating discrete choice models, because it both provides a theoretically strong foundation and facilitates intuitions about the model, which in turn makes it easy to consider various sorts of extensions. (See the example below.)

The choice of the type-1 extreme value distribution seems fairly arbitrary, but it makes the mathematics work out, and it may be possible to justify its use through rational choice theory.

It turns out that this model is equivalent to the previous model, although this seems non-obvious, since there are now two sets of regression coefficients and error variables, and the error variables have a different distribution. In fact, this model reduces directly to the previous one with the following substitutions:

$$\beta = \beta_1 - \beta_0$$

$$\varepsilon = \varepsilon_1 - \varepsilon_0$$

An intuition for this comes from the fact that, since we choose based on the maximum of two values, only their difference matters, not the exact values — and this effectively removes one degree of freedom. Another critical fact is that the difference of two type-1 extreme-value-distributed variables is a logistic distribution, i.e. $\varepsilon = \varepsilon_1 - \varepsilon_0 \sim \mathbf{Logistic}(0, 1)$. We can demonstrate the equivalent as follows:

$$\begin{aligned} \Pr(Y_i = 1 \mid \mathbf{X}_i) &= \Pr\left(Y_i^{1*} > Y_i^{0*} \mid \mathbf{X}_i\right) \\ &= \Pr\left(Y_i^{1*} - Y_i^{0*} > 0 \mid \mathbf{X}_i\right) \\ &= \Pr\left(\beta_1 \cdot \mathbf{X}_i + \varepsilon_1 - (\beta_0 \cdot \mathbf{X}_i + \varepsilon_0) > 0\right) \\ &= \Pr\left((\beta_1 \cdot \mathbf{X}_i - \beta_0 \cdot \mathbf{X}_i) + (\varepsilon_1 - \varepsilon_0) > 0\right) \\ &= \Pr\left((\beta_1 - \beta_0) \cdot \mathbf{X}_i + (\varepsilon_1 - \varepsilon_0) > 0\right) \\ &= \Pr\left((\beta_1 - \beta_0) \cdot \mathbf{X}_i + \varepsilon > 0\right) && \text{(substitute } \varepsilon \text{ as above)} \\ &= \Pr\left(\beta \cdot \mathbf{X}_i + \varepsilon > 0\right) && \text{(substitute } \beta \text{ as above)} \\ &= \Pr\left(\varepsilon > -\beta \cdot \mathbf{X}_i\right) && \text{(now, same as above model)} \\ &= \Pr\left(\varepsilon < \beta \cdot \mathbf{X}_i\right) \\ &= \text{logit}^{-1}(\beta \cdot \mathbf{X}_i) \\ &= p_i \end{aligned}$$

Example

As an example, consider a province-level election where the choice is between a right-of-center party, a left-of-center party, and a secessionist party (e.g. the Parti Québécois, which wants Quebec to secede from Canada). We would then use three latent variables, one for each choice. Then, in accordance with utility theory, we can then interpret the latent variables as expressing the utility that results from making each of the choices. We can also interpret the regression coefficients as indicating the strength that the associated factor (i.e. explanatory variable) has in contributing to the utility — or more correctly, the amount by which a unit change in an explanatory variable changes the utility of a given choice. A voter might expect that the right-of-center party would lower taxes, especially on rich people. This would give low-income people no benefit, i.e. no change in utility (since they usually don't pay taxes); would cause moderate benefit (i.e. somewhat more money, or moderate utility increase) for middle-incoming people; would cause significant benefits for high-income people. On the other hand, the left-of-center party might be expected to raise taxes and offset it with increased welfare and other assistance for the lower and middle classes. This would cause significant positive benefit to low-income people, perhaps weak benefit to middle-income people, and significant negative benefit to high-income people. Finally, the secessionist party would take no direct actions on the economy, but simply secede. A low-income or middle-income voter might expect basically no clear utility gain or loss from this, but a high-income voter might expect negative utility since he/she is likely to own companies, which will have a harder time doing business in such an environment and probably lose money.

These intuitions can be expressed as follows:

Estimated strength of regression coefficient for different outcomes (party choices) and different values of explanatory variables

	Center-right	Center-left	Secessionist
High-income	strong +	strong −	strong −
Middle-income	moderate +	weak +	none
Low-income	none	strong +	none

This clearly shows that

- Separate sets of regression coefficients need to exist for each choice. When phrased in terms of utility, this can be seen very easily. Different choices have different effects on net utility; furthermore, the effects vary in complex ways that depend on the characteristics of each individual, so there need to be separate sets of coefficients for each characteristic, not simply a single extra per-choice characteristic.
- Even though income is a continuous variable, its effect on utility is too complex for it to be treated as a single variable. Either it needs to be directly split up into ranges, or higher powers of income need to be added so that polynomial regression on income is effectively done.

As a "log-linear" model

Yet another formulation combines the two-way latent variable formulation above with the original formulation higher up without latent variables, and in the process provides a link to one of the standard formulations of the multinomial logit.

Here, instead of writing the logit of the probabilities p_i as a linear predictor, we separate the linear predictor into two, one for each of the two outcomes:

$$\begin{aligned}\ln \Pr(Y_i = 0) &= \boldsymbol{\beta}_0 \cdot \mathbf{X}_i - \ln Z \\ \ln \Pr(Y_i = 1) &= \boldsymbol{\beta}_1 \cdot \mathbf{X}_i - \ln Z\end{aligned}$$

Note that two separate sets of regression coefficients have been introduced, just as in the two-way latent variable model, and the two equations appear a form that writes the logarithm of the associated probability as a linear predictor, with an extra term $-\ln Z$ at the end. This term, as it turns out, serves as the normalizing factor ensuring that the result is a distribution. This can be seen by exponentiating both sides:

$$\begin{aligned}\Pr(Y_i = 0) &= \frac{1}{Z} e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} \\ \Pr(Y_i = 1) &= \frac{1}{Z} e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}\end{aligned}$$

In this form it is clear that the purpose of Z is to ensure that the resulting distribution over Y_i is in fact a probability distribution, i.e. it sums to 1. This means that Z is simply the sum of all un-normalized probabilities, and by dividing each probability by Z , the probabilities become "normalized". That is:

$$Z = e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}$$

and the resulting equations are

$$\begin{aligned}\Pr(Y_i = 0) &= \frac{e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i}}{e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}} \\ \Pr(Y_i = 1) &= \frac{e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}}{e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}}.\end{aligned}$$

Or generally:

$$\Pr(Y_i = c) = \frac{e^{\boldsymbol{\beta}_c \cdot \mathbf{X}_i}}{\sum_h e^{\boldsymbol{\beta}_h \cdot \mathbf{X}_i}}$$

This shows clearly how to generalize this formulation to more than two outcomes, as in multinomial logit. Note that this general formulation is exactly the Softmax function as in

$$\Pr(Y_i = c) = \text{softmax}(c, \boldsymbol{\beta}_0 \cdot \mathbf{X}_i, \boldsymbol{\beta}_1 \cdot \mathbf{X}_i, \dots).$$

In order to prove that this is equivalent to the previous model, note that the above model is overspecified, in that $\Pr(Y_i = 0)$ and $\Pr(Y_i = 1)$ cannot be independently specified: rather $\Pr(Y_i = 0) + \Pr(Y_i = 1) = 1$ so knowing one automatically determines the other. As a result, the model is nonidentifiable, in that multiple combinations of $\boldsymbol{\beta}_0$ and $\boldsymbol{\beta}_1$ will produce the same probabilities for all possible explanatory variables. In fact, it can be seen that adding any constant vector to both of them will produce the same probabilities:

$$\begin{aligned}\Pr(Y_i = 1) &= \frac{e^{(\boldsymbol{\beta}_1 + \mathbf{C}) \cdot \mathbf{X}_i}}{e^{(\boldsymbol{\beta}_0 + \mathbf{C}) \cdot \mathbf{X}_i} + e^{(\boldsymbol{\beta}_1 + \mathbf{C}) \cdot \mathbf{X}_i}} \\ &= \frac{e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i} e^{\mathbf{C} \cdot \mathbf{X}_i}}{e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} e^{\mathbf{C} \cdot \mathbf{X}_i} + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i} e^{\mathbf{C} \cdot \mathbf{X}_i}} \\ &= \frac{e^{\mathbf{C} \cdot \mathbf{X}_i} e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}}{e^{\mathbf{C} \cdot \mathbf{X}_i} (e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i})} \\ &= \frac{e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}}{e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}}.\end{aligned}$$

As a result, we can simplify matters, and restore identifiability, by picking an arbitrary value for one of the two vectors. We choose to set $\boldsymbol{\beta}_0 = \mathbf{0}$. Then,

$$e^{\boldsymbol{\beta}_0 \cdot \mathbf{X}_i} = e^{\mathbf{0} \cdot \mathbf{X}_i} = 1$$

and so

$$\Pr(Y_i = 1) = \frac{e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}}{1 + e^{\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}} = \frac{1}{1 + e^{-\boldsymbol{\beta}_1 \cdot \mathbf{X}_i}} = p_i$$

which shows that this formulation is indeed equivalent to the previous formulation. (As in the two-way latent variable formulation, any settings where $\boldsymbol{\beta} = \boldsymbol{\beta}_1 - \boldsymbol{\beta}_0$ will produce equivalent results.)

Note that most treatments of the multinomial logit model start out either by extending the "log-linear" formulation presented here or the two-way latent variable formulation presented above, since both clearly show the way that the model could be extended to multi-way outcomes. In general, the presentation with latent variables is more common in econometrics and political science, where discrete choice models and utility theory reign, while the "log-linear" formulation here is more common in computer science, e.g. machine learning and natural language processing.

As a single-layer perceptron

The model has an equivalent formulation

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i})}}.$$

This functional form is commonly called a single-layer perceptron or single-layer artificial neural network. A single-layer neural network computes a continuous output instead of a step function. The derivative of p_i with respect to $X = (x_1, ..., x_k)$ is computed from the general form:

$$y = \frac{1}{1 + e^{-f(X)}}$$

where $f(X)$ is an analytic function in X . With this choice, the single-layer neural network is identical to the logistic regression model. This function has a continuous derivative, which allows it to be used in backpropagation. This function is also preferred because its derivative is easily calculated:

$$\frac{dy}{dX} = y(1 - y) \frac{df}{dX}.$$

In terms of binomial data

A closely related model assumes that each i is associated not with a single Bernoulli trial but with n_i independent identically distributed trials, where the observation Y_i is the number of successes observed (the sum of the individual Bernoulli-distributed random variables), and hence follows a binomial distribution:

$$Y_i \sim \text{Bin}(n_i, p_i), \text{ for } i = 1, \dots, n$$

An example of this distribution is the fraction of seeds (p_i) that germinate after n_i are planted.

In terms of expected values, this model is expressed as follows:

$$p_i = \text{E} \left[\frac{Y_i}{n_i} \mid \mathbf{X}_i \right],$$

so that

$$\text{logit} \left(\text{E} \left[\frac{Y_i}{n_i} \mid \mathbf{X}_i \right] \right) = \text{logit}(p_i) = \ln \left(\frac{p_i}{1 - p_i} \right) = \boldsymbol{\beta} \cdot \mathbf{X}_i,$$

Or equivalently:

$$\Pr(Y_i = y \mid \mathbf{X}_i) = \binom{n_i}{y} p_i^y (1 - p_i)^{n_i - y} = \binom{n_i}{y} \left(\frac{1}{1 + e^{-\boldsymbol{\beta} \cdot \mathbf{X}_i}} \right)^y \left(1 - \frac{1}{1 + e^{-\boldsymbol{\beta} \cdot \mathbf{X}_i}} \right)^{n_i - y}.$$

This model can be fit using the same sorts of methods as the above more basic model.

Bayesian

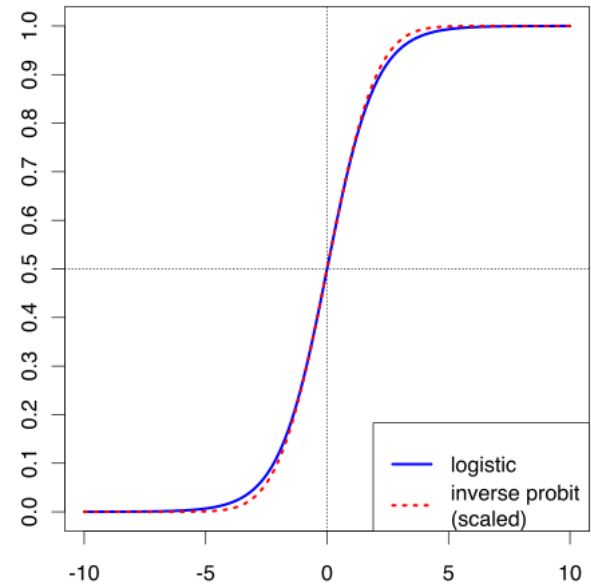
In a Bayesian statistics context, prior distributions are normally placed on the regression coefficients, usually in the form of Gaussian distributions. There is no conjugate prior of the likelihood function in logistic regression. When Bayesian inference was performed analytically, this made the posterior distribution difficult to calculate except in very low dimensions. Now, though, automatic software such as OpenBUGS, JAGS, PyMC3 or Stan allows these posteriors to be computed using simulation, so lack of conjugacy is not a concern. However, when the sample size or the number of parameters is large, full Bayesian simulation can be slow, and people often use approximate methods such as variational Bayesian methods and expectation propagation.

History

A detailed history of the logistic regression is given in Cramer (2002). The logistic function was developed as a model of population growth and named "logistic" by Pierre Franois Verhulst in the 1830s and 1840s, under the guidance of Adolphe Quetelet; see Logistic function § History for details.^[35] In his earliest paper (1838), Verhulst did not specify how he fit the curves to the data.^{[36][37]} In his more detailed paper (1845), Verhulst determined the three parameters of the model by making the curve pass through three observed points, which yielded poor predictions.^{[38][39]}

The logistic function was independently developed in chemistry as a model of autocatalysis (Wilhelm Ostwald, 1883).^[40] An autocatalytic reaction is one in which one of the products is itself a catalyst for the same reaction, while the supply of one of the reactants is fixed. This naturally gives rise to the logistic equation for the same reason as population growth: the reaction is self-reinforcing but constrained.

The logistic function was independently rediscovered as a model of population growth in 1920 by Raymond Pearl and Lowell Reed, published as Pearl & Reed (1920), which led to its use in modern statistics. They were initially unaware of Verhulst's work and presumably learned about it from L. Gustave du Pasquier, but they gave him little credit and did not adopt his terminology.^[41] Verhulst's priority was acknowledged and the term "logistic" revived by Udny Yule in 1925 and has been followed since.^[42] Pearl and Reed first applied the model to the population of the United States, and also initially fitted the curve by making it pass through three points; as with Verhulst, this again yielded poor results.^[43]



Comparison of logistic function with a scaled inverse probit function (i.e. the CDF of the normal distribution), comparing $\sigma(\boldsymbol{x})$ vs. $\Phi(\sqrt{\frac{\pi}{8}} \boldsymbol{x})$, which makes the slopes the same at the origin. This shows the heavier tails of the logistic distribution.

In the 1930s, the [probit model](#) was developed and systematized by [Chester Ittner Bliss](#), who coined the term "probit" in [Bliss \(1934\)](#), and by [John Gaddum](#) in [Gaddum \(1933\)](#), and the model fit by [maximum likelihood estimation](#) by [Ronald A. Fisher](#) in [Fisher \(1935\)](#), as an addendum to Bliss's work. The probit model was principally used in [bioassay](#), and had been preceded by earlier work dating to 1860; see [Probit model § History](#). The probit model influenced the subsequent development of the logit model and these models competed with each other.^[44]

The logistic model was likely first used as an alternative to the probit model in bioassay by [Edwin Bidwell Wilson](#) and his student [Jane Worcester](#) in [Wilson & Worcester \(1943\)](#).^[45] However, the development of the logistic model as a general alternative to the probit model was principally due to the work of [Joseph Berkson](#) over many decades, beginning in [Berkson \(1944\)](#), where he coined "logit", by analogy with "probit", and continuing through [Berkson \(1951\)](#) and following years.^[46] The logit model was initially dismissed as inferior to the probit model, but "gradually achieved an equal footing with the logit",^[47] particularly between 1960 and 1970. By 1970, the logit model achieved parity with the probit model in use in statistics journals and thereafter was surpassed it. This relative popularity was due to the adoption of the logit outside of bioassay, rather than displacing the probit within bioassay, and its informal use in practice; the logit's popularity is credited to the logit model's computational simplicity, mathematical properties, and generality, allowing its use in varied fields.^[48]

Various refinements occurred during that time, notably by [David Cox](#), as in [Cox \(1958\)](#).^[1]

The multinomial logit model was introduced independently in [Cox \(1966\)](#) and [Thiel \(1969\)](#), which greatly increased the scope of application and the popularity of the logit model.^[49] In 1973 [Daniel McFadden](#) linked the multinomial logit to the theory of [discrete choice](#), specifically [Luce's choice axiom](#), showing that the multinomial logit followed from the assumption of [independence of irrelevant alternatives](#) and interpreting odds of alternatives as relative preferences;^[50] this gave a theoretical foundation for the logistic regression.^[49]

Extensions

There are large numbers of extensions:

- [Multinomial logistic regression](#) (or **multinomial logit**) handles the case of a multi-way [categorical](#) dependent variable (with unordered values, also called "classification"). Note that the general case of having dependent variables with more than two values is termed *polytomous regression*.
- [Ordered logistic regression](#) (or **ordered logit**) handles [ordinal](#) dependent variables (ordered values).
- [Mixed logit](#) is an extension of multinomial logit that allows for correlations among the choices of the dependent variable.
- An extension of the logistic model to sets of interdependent variables is the [conditional random field](#).
- [Conditional logistic regression](#) handles [matched](#) or [stratified](#) data when the strata are small. It is mostly used in the analysis of [observational studies](#).

Software

Most [statistical software](#) can do binary logistic regression.

- [SPSS](#)
 - [\[2\]](#) (<http://www-01.ibm.com/support/docview.wss?uid=swg21475013>) for basic logistic regression.
- [Stata](#)
- [SAS](#)
 - [PROC LOGISTIC](#) (https://support.sas.com/documentation/cdl/en/statug/63347/HTML/default/viewer.htm#logistic_toc.htm) for basic logistic regression.
 - [PROC CATMOD](#) (https://support.sas.com/documentation/cdl/en/statug/63347/HTML/default/viewer.htm#statug_catmod_sect003.htm) when all the variables are categorical.
 - [PROC GLIMMIX](#) (https://support.sas.com/documentation/cdl/en/statug/63033/HTML/default/viewer.htm#glimmix_toc.htm) for [multilevel model](#) logistic regression.
- [R](#)
 - [glm](#) in the [stats](#) package (using family = binomial)^[51]
 - [lrm](#) in the [rms](#) package (<https://cran.r-project.org/web/packages/rms>)
 - [GLMNET](#) package for an efficient implementation regularized logistic regression
 - [lmer](#) for mixed effects logistic regression
 - [Rfast](#) package command [gm_logistic](#) for fast and heavy calculations involving large scale data.
 - [arm](#) package for bayesian logistic regression
- [Python](#)
 - [Logit](#) (http://www.statsmodels.org/dev/generated/statsmodels.discrete.discrete_model.Logit.html) in the [Statsmodels](#) module.
 - [LogisticRegression](#) (http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html) in the [Scikit-learn](#) module.
 - [LogisticRegressor](#) (https://www.tensorflow.org/versions/master/api_docs/python/tf/contrib/learn/LogisticRegressor) in the [TensorFlow](#) module.
 - Full example of logistic regression in the Theano tutorial [\[3\]](#) (<http://deeplearning.net/software/theano/tutorial/examples.html>)
 - Bayesian Logistic Regression with ARD prior [code](#) (https://github.com/AmazaspShumik/sklearn-bayes/blob/master/skbayes/rvm_ard_models/fast_rvm.py), [tutorial](#) (http://github.com/AmazaspShumik/sklearn-bayes/blob/master/ipython_notebooks_tutorials/rvm_ard/ard_classification_demo.ipynb)
 - Variational Bayes Logistic Regression with ARD prior [code](#) (https://github.com/AmazaspShumik/sklearn-bayes/blob/master/skbayes/rvm_ard_models/vrvm.py) , [tutorial](#) (https://github.com/AmazaspShumik/sklearn-bayes/blob/master/ipython_notebooks_tutorials/rvm_ard/vbard_classification.ipynb)
 - Bayesian Logistic Regression [code](#) (https://github.com/AmazaspShumik/sklearn-bayes/blob/master/skbayes/linear_models/bayes_logistic.py), [tutorial](#) (https://github.com/AmazaspShumik/sklearn-bayes/blob/master/ipython_notebooks_tutorials/linear_models/bayesian_logistic_regression_demo.ipynb)
- [NCSS](#)
 - [Logistic Regression in NCSS](#) (http://ncss.wpengine.netdna-cdn.com/wp-content/themes/ncss/pdf/Procedures/NCSS/Logistic_Regression.pdf)
- [Matlab](#)
 - [mnrfit](#) in the [Statistics and Machine Learning Toolbox](#) (with "incorrect" coded as 2 instead of 0)
- [Java](#) ([JVM](#))
 - [LibLinear](#)
 - [Apache Flink](#) (https://ci.apache.org/projects/flink/flink-docs-release-1.7/dev/libs/ml/multiple_linear_regression.html)
 - [Apache Spark](#)
 - [SparkML](#) supports Logistic Regression

Notably, [Microsoft Excel](#)'s statistics extension package does not include it.

See also

- Logistic function
- Discrete choice
- Jarrow–Turnbull model
- Limited dependent variable
- Multinomial logit model
- Ordered logit
- Hosmer–Lemeshow test
- Brier score
- mlpack - contains a C++ implementation of logistic regression
- Local case-control sampling
- Logistic model tree

References

1.

Walker, SH; Duncan, DB (1967). "Estimation of the probability of an event as a function of several independent variables". *Biometrika*. **54** (1/2): 167–178. doi:10.2307/2333860 (https://doi.org/10.2307%2F2333860). JSTOR 2333860 (http s://www.jstor.org/stable/2333860).

2.

Cramer 2002, p. 8.

3.

Boyd, C. R.; Tolson, M. A.; Copes, W. S. (1987). "Evaluating trauma care: The TRISS method. Trauma Score and the Injury Severity Score". *The Journal of Trauma*. **27** (4): 370–378. doi:10.1097/00005373-198704000-00005 (https://doi.org/10.1097%2F000 05373-198704000-00005). PMID 3106646 (https://www.ncbi.nlm.nih.gov/pubmed/ 3106646).

4.

Kologlu, M.; Elker, D.; Altun, H.; Sayek, I. (2001). "Validation of MPI and PIA II in two different groups of patients with secondary peritonitis". *Hepato-Gastroenterology*. **48** (37): 147–51. PMID 11268952 (https://www.ncbi.nlm.nih.gov/pubmed/1126895 2).

5.

Biondo, S.; Ramos, E.; Deiros, M.; Ragué, J. M.; De Oca, J.; Moreno, P.; Farran, L.; Jaurieta, E. (2000). "Prognostic factors for mortality in left colonic peritonitis: A new scoring system". *Journal of the American College of Surgeons*. **191** (6): 635–42. doi:10.1016/S1072-7515(00)00758-4 (https://doi.org/10.1016%2FS1072-7515%280 0%2900758-4). PMID 11129812 (https://www.ncbi.nlm.nih.gov/pubmed/11129812).

6.

Marshall, J. C.; Cook, D. J.; Christou, N. V.; Bernard, G. R.; Sprung, C. L.; Sibbald, W. J. (1995). "Multiple organ dysfunction score: A reliable descriptor of a complex clinical outcome". *Critical Care Medicine*. **23** (10): 1638–52. doi:10.1097/00003246-199510000-00007 (https://doi.org/10.1097%2F00003246-199510000-00007). PMID 7587228 (https://www.ncbi.nlm.nih.gov/pubmed/7587228).

7.

Le Gall, J. R.; Lemeshow, S.; Saulnier, F. (1993). "A new Simplified Acute Physiology Score (SAPS II) based on a European/North American multicenter study". *JAMA*. **270** (24): 2957–63. doi:10.1001/jama.1993.03510240069035 (https://doi.org/10.1001%2 Fjama.1993.03510240069035). PMID 8254858 (https://www.ncbi.nlm.nih.gov/pub med/8254858).

8.

David A. Freedman (2009). *Statistical Models: Theory and Practice*. Cambridge University Press. p. 128.

9.

Truett, J; Cornfield, J; Kannel, W (1967). "A multivariate analysis of the risk of coronary heart disease in Framingham". *Journal of Chronic Diseases*. **20** (7): 511–24. doi:10.1016/0021-9681(67)90082-3 (https://doi.org/10.1016%2F0021-9681%2867% 2990082-3). PMID 6028270 (https://www.ncbi.nlm.nih.gov/pubmed/6028270).

10.

Harrell, Frank E. (2001). *Regression Modeling Strategies* (2nd ed.). Springer-Verlag. ISBN 978-0-387-95232-1.

11.

M. Strano; B.M. Colosimo (2006). "Logistic regression analysis for experimental determination of forming limit diagrams". *International Journal of Machine Tools and Manufacture*. **46** (6): 673–682. doi:10.1016/j.ijmachtools.2005.07.005 (https://do i.org/10.1016%2Fj.ijmachtools.2005.07.005).

12.

Palei, S. K.; Das, S. K. (2009). "Logistic regression model for prediction of roof fall risks in bord and pillar workings in coal mines: An approach". *Safety Science*. **47**: 88–96. doi:10.1016/j.ssci.2008.01.002 (https://doi.org/10.1016%2Fj.ssci.2008.01.00 2).

13.

Berry, Michael J.A (1997). *Data Mining Techniques For Marketing, Sales and Customer Support*. Wiley. p. 10.

14.

Hosmer, David W.; Lemeshow, Stanley (2000). *Applied Logistic Regression* (2nd ed.). Wiley. ISBN 978-0-471-35632-5.

15.

Harrell, Frank E. (2015). *Regression Modeling Strategies*. Springer Series in Statistics (2nd ed.). New York; Springer. doi:10.1007/978-3-319-19425-7 (https://doi.org/10.1 007%2F978-3-319-19425-7). ISBN 978-3-319-19424-0.

16.

Rodríguez, G. (2007). *Lecture Notes on Generalized Linear Models*. pp. Chapter 3, page 45 – via http://data.princeton.edu/wws509/notes/.

17.

Gareth James; Daniela Witten; Trevor Hastie; Robert Tibshirani (2013). *An Introduction to Statistical Learning* (http://www-bcf.usc.edu/~gareth/ISL/). Springer. p. 6.

18.

Pohar, Maja; Blas, Mateja; Turk, Sandra (2004). "Comparison of Logistic Regression and Linear Discriminant Analysis: A Simulation Study" (https://www.researchgate.n et/publication/229021894). *Metodološki Zvezki*. **1** (1).

19.

"How to Interpret Odds Ratio in Logistic Regression?" (https://stats.idre.ucla.edu/st ata/faq/how-do-i-interpret-odds-ratios-in-logistic-regression/). Institute for Digital Research and Education.

20.

Everitt, Brian (1998). *The Cambridge Dictionary of Statistics*. Cambridge, UK New York: Cambridge University Press. ISBN 978-0521593465.

21.

Ng, Andrew (2000). "CS229 Lecture Notes" (http://akademik.bahcesehir.edu.tr/~tev fik/courses/cmp5101/cs229-notes1.pdf) (PDF). *CS229 Lecture Notes*: 16–19.

22.

Van Smeden, M.; De Groot, J. A.; Moons, K. G.; Collins, G. S.; Altman, D. G.; Eijkemans, M. J.; Reitsma, J. B. (2016). "No rationale for 1 variable per 10 events criterion for binary logistic regression analysis" (https://www.ncbi.nlm.nih.gov/pm c/articles/PMC5122171). *BMC Medical Research Methodology*. **16** (1): 163. doi:10.1186/s12874-016-0267-3 (https://doi.org/10.1186%2Fs12874-016-0267-3). PMC 5122171 (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5122171). PMID 27881078 (https://www.ncbi.nlm.nih.gov/pubmed/27881078).

23.

Peduzzi, P; Concato, J; Kemper, E; Holford, TR; Feinstein, AR (December 1996). "A simulation study of the number of events per variable in logistic regression analysis". *Journal of Clinical Epidemiology*. **49** (12): 1373–9. doi:10.1016/s0895-4356(96)00236-3 (https://doi.org/10.1016%2Fs0895-4356%2896%2900236-3). PMID 8970487 (https://www.ncbi.nlm.nih.gov/pubmed/8970487).

24.

Vittinghoff, E.; McCulloch, C. E. (12 January 2007). "Relaxing the Rule of Ten Events per Variable in Logistic and Cox Regression" (http://aje.oxfordjournals.org/conten t/165/6/710.full). *American Journal of Epidemiology*. **165** (6): 710–718. doi:10.1093/aje/kwk052 (https://doi.org/10.1093%2Faje%2Fkwk052). PMID 17182981 (https://www.ncbi.nlm.nih.gov/pubmed/17182981).

25.

van der Ploeg, Tjeerd; Austin, Peter C.; Steyerberg, Ewout W. (2014). "Modern modelling techniques are data hungry: a simulation study for predicting dichotomous endpoints" (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC428955 3). *BMC Medical Research Methodology*. **14**: 137. doi:10.1186/1471-2288-14-137 (http s://doi.org/10.1186%2F1471-2288-14-137). PMC 4289553 (https://www.ncbi.nlm.ni h.gov/pmc/articles/PMC4289553). PMID 25532820 (https://www.ncbi.nlm.nih.gov/ pubmed/25532820).

26.

Menard, Scott W. (2002). *Applied Logistic Regression* (2nd ed.). SAGE. ISBN 978-0-7619-2208-7.

27.

Murphy, Kevin P. (2012). *Machine Learning – A Probabilistic Perspective*. The MIT Press. pp. 245pp. ISBN 978-0-262-01802-9.

28.

Greene, William N. (2003). *Econometric Analysis* (Fifth ed.). Prentice-Hall. ISBN 978-0-13-066189-0.

29.

Cohen, Jacob; Cohen, Patricia; West, Steven G.; Aiken, Leona S. (2002). *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences* (3rd ed.). Routledge. ISBN 978-0-8058-2223-6.

30.

Allison, Paul D. "Measures of Fit for Logistic Regression" (https://support.sas.com/r esources/papers/proceedings14/1485-2014.pdf) (PDF). Statistical Horizons LLC and the University of Pennsylvania.

31.

Tjur, Tue (2009). "Coefficients of determination in logistic regression models". *American Statistician*: 366–372. doi:10.1198/tast.2009.08210 (https://doi.org/10.119 8%2Ftast.2009.08210).

32.

Hosmer, D.W. (1997). "A comparison of goodness-of-fit tests for the logistic regression model". *Stat Med*. **16** (9): 965–980. doi:10.1002/(sici)1097-0258(19970515)16:9<965::aid-sim509>3.3.co;2-f (https://doi.org/10.1002%2F%28sic i%291097-0258%2819970515%2916%3A9%3C965%3A%3Aaid-sim509%3E3.3.co%3 B2-f).

33.

https://class.stanford.edu/c4x/HumanitiesScience/StatLearning/asset/classification.p slide 16

34.

Malouf, Robert (2002). "A comparison of algorithms for maximum entropy parameter estimation" (https://dl.acm.org/citation.cfm?id=1118871). *Proceedings of the Sixth Conference on Natural Language Learning (CoNLL-2002)*. pp. 49–55. doi:10.3115/1118853.1118871 (https://doi.org/10.3115%2F1118853.1118871).

35.

Cramer 2002, pp. 3–5.

36.

Verhulst, Pierre-François (1838). "Notice sur la loi que la population poursuit dans son accroissement" (https://books.google.com/?id=8GsEAAAAYAAJ) (PDF). *Correspondance Mathématique et Physique*. **10**: 113–121. Retrieved 3 December 2014.

37.

Cramer 2002, p. 4, "He did not say how he fitted the curves."

38. Verhulst, Pierre-François (1845). "Recherches mathématiques sur la loi d'accroissement de la population" (http://gdz.sub.uni-goettingen.de/dms/load/img/?PPN=PPN129323640_0018&DMDID=dmdlog7) [Mathematical Researches into the Law of Population Growth Increase]. *Nouveaux Mémoires de l'Académie Royale des Sciences et Belles-Lettres de Bruxelles*. **18**. Retrieved 2013-02-18.

39. Cramer 2002, p. 4.

40. Cramer 2002, p. 7.

41. Cramer 2002, p. 6.

42. Cramer 2002, p. 6–7.

43. Cramer 2002, p. 5.

44. Cramer 2002, p. 7–9.

45. Cramer 2002, p. 9.

46. Cramer 2002, p. 8, "As far as I can see the introduction of the logistics as an alternative to the normal probability function is the work of a single person, Joseph Berkson (1899–1982), ..."

47. Cramer 2002, p. 11.

48. Cramer 2002, p. 10–11.

49. Cramer, p. 13.

50. McFadden, Daniel (1973). "Conditional Logit Analysis of Qualitative Choice Behavior" (https://web.archive.org/web/20181127110612/https://eml.berkeley.edu/reprints/mcfadden/zarembka.pdf) (PDF). In P. Zarembka (ed.). *Frontiers in Econometrics*. New York: Academic Press. pp. 105–142. Archived from the original (https://eml.berkeley.edu/reprints/mcfadden/zarembka.pdf) (PDF) on 2018-11-27. Retrieved 2019-04-20.

51. Gelman, Andrew; Hill, Jennifer (2007). *Data Analysis Using Regression and Multilevel/Hierarchical Models* (https://books.google.com/books?id=IV3DIIdV0F9AC&pg=PA79). New York: Cambridge University Press. pp. 79–108. ISBN 978-0-521-68689-1.

Further reading

■ Cox, David R. (1958). "The regression analysis of binary sequences (with discussion)". *J Roy Stat Soc B*. **20** (2): 215–242. JSTOR 2983890 (https://www.jstor.org/stable/2983890).

■ Cox, David R. (1966). "Some procedures connected with the logistic qualitative response curve". In F. N. David (1966) (ed.). *Research Papers in Probability and Statistics (Festschrift for J. Neyman)*. London: Wiley. pp. 55–71.

■ Cramer, J. S. (2002). *The origins of logistic regression* (https://papers.tinbergen.nl/02119.pdf) (PDF) (Technical report). **119**. Tinbergen Institute. pp. 167–178. doi:10.2139/ssrn.360300 (https://doi.org/10.2139%2Fssrn.360300).

■ Published in: Cramer, J. S. (2004). "The early origins of the logit model". *Studies in History and Philosophy of Science Part C: Studies in History and Philosophy of Biological and Biomedical Sciences*. **35** (4): 613–626. doi:10.1016/j.shpsc.2004.09.003 (https://doi.org/10.1016%2Fj.shpsc.2004.09.003).

■ Thiel, Henri (1969). "A Multinomial Extension of the Linear Logit Model". *International Economic Review*. **10** (3): 251–59. doi:10.2307/2525642 (https://doi.org/10.2307%2F2525642). JSTOR 2525642 (https://www.jstor.org/stable/2525642).

■ Wilson, E.B.; Worcester, J. (1943). "The Determination of L.D.50 and Its Sampling Error in Bio-Assay" (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1078563). *Proceedings of the National Academy of Sciences of the United States of America*. **29** (2): 79–85. Bibcode:1943PNAS...29...79W (http://adsabs.harvard.edu/abs/1943PNA S...29...79W). doi:10.1073/pnas.29.2.79 (https://doi.org/10.1073%2Fpnas.29.2.79). PMC 1078563 (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1078563). PMID 16588606 (https://www.ncbi.nlm.nih.gov/pubmed/16588606).

■ Agresti, Alan. (2002). *Categorical Data Analysis*. New York: Wiley-Interscience. ISBN 978-0-471-36093-3.

■ Amemiya, Takeshi (1985). "Qualitative Response Models" (https://books.google.co m/books?id=0bzGQE14CwEC&pg=PA267). *Advanced Econometrics*. Oxford: Basil Blackwell. pp. 267–359. ISBN 978-0-631-13345-2.

■ Balakrishnan, N. (1991). *Handbook of the Logistic Distribution*. Marcel Dekker, Inc. ISBN 978-0-8247-8587-1.

■ Gouriéroux, Christian (2000). "The Simple Dichotomy" (https://books.google.com/b ooks?id=dE2prs_U0QMC&pg=PA6). *Econometrics of Qualitative Dependent Variables*. New York: Cambridge University Press. pp. 6–37. ISBN 978-0-521-58985-7.

■ Greene, William H. (2003). *Econometric Analysis, fifth edition*. Prentice Hall. ISBN 978-0-13-066189-0.

■ Hilbe, Joseph M. (2009). *Logistic Regression Models*. Chapman & Hall/CRC Press. ISBN 978-1-4200-7575-5.

■ Hosmer, David (2013). *Applied logistic regression*. Hoboken, New Jersey: Wiley. ISBN 978-0470582473.

■ Howell, David C. (2010). *Statistical Methods for Psychology, 7th ed.* Belmont, CA; Thomson Wadsworth. ISBN 978-0-495-59786-5.

■ Peduzzi, P.; J. Concato; E. Kemper; T.R. Holford; A.R. Feinstein (1996). "A simulation study of the number of events per variable in logistic regression analysis". *Journal of Clinical Epidemiology*. **49** (12): 1373–1379. doi:10.1016/s0895-4356(96)00236-3 (ht tps://doi.org/10.1016%2Fs0895-4356%2896%2900236-3). PMID 8970487 (https://w ww.ncbi.nlm.nih.gov/pubmed/8970487).

■ Berry, Michael J.A.; Linoff, Gordon (1997). *Data Mining Techniques For Marketing, Sales and Customer Support*. Wiley.

External links

- Econometrics Lecture (topic: Logit model) (https://www.youtube.com/watch?v=JvioZoK1f4o&t=64m48s) on YouTube by Mark Thoma

■ Logistic Regression tutorial (http://www.omidrouhani.com/research/logisticregression/html/logisticregression.htm)

■ mlelr (https://czep.net/stat/mlelr.html): software in C for teaching purposes

Retrieved from "https://en.wikipedia.org/w/index.php?title=Logistic_regression&oldid=899297078"

This page was last edited on 29 May 2019, at 03:41 (UTC).

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the [Terms of Use](#) and [Privacy Policy](#). Wikipedia® is a registered trademark of the [Wikimedia Foundation, Inc.](#), a non-profit organization.