

## Assumptions of logistic regression

- Logistic regression does not assume a linear relationship between the dependent and independent variables.
- The dependent variable must be a dichotomy (2 categories). (*Remark: Dichotomous refers to two outcomes. Multichotomous refers to more than two outcomes.*)
- Examples of dichotomous variables include sex (two groups: “males” and “females”), presence of heart disease (two groups: “yes” and “no”), personality type (two groups: “introversion” or “extroversion”), body composition (two groups: “obese” or “not obese”), and so forth.
- **Important:** The independent variables need not be interval, nor normally distributed, nor linearly related, nor of equal variance within each group.
- The categories (groups) must be mutually exclusive and exhaustive; a case can only be in one group and every case must be a member of one of the groups.
- Larger samples are needed than for linear regression because maximum likelihood coefficients are large sample estimates. A minimum of 50 cases per predictor is recommended.

## Maximum Likelihood Estimation

- Maximum likelihood estimation, MLE, is the method used to calculate the logit coefficients. This contrasts to the use of ordinary least squares (OLS) estimation of coefficients in regression. OLS seeks to minimize the sum of squared distances of the data points to the regression line.
- MLE seeks to maximize the log likelihood, LL, which reflects how likely it is (the odds) that the observed values of the dependent may be predicted from the observed values of the independents. (Equivalently MLE seeks to minimize the -2LL value.)
- MLE is an iterative algorithm which starts with an initial arbitrary “guesstimate” of what the logit coefficients should be, the MLE algorithm determines the direction and size change in the logit coefficients which will increase LL.
- After this initial function is estimated, the residuals are tested and a re-estimate is made with an improved function, and the process is repeated (usually about a half-dozen times) until convergence is reached (that is, until LL does not change significantly). There are several alternative convergence criteria.

## Log Likelihood

- A “*likelihood*” is a probability, specifically the probability that the observed values of the dependent may be predicted from the observed values of the independents.
- Like any probability, the likelihood varies from 0 to 1. The log likelihood (LL) is its log and varies from 0 to minus infinity (it is negative because the log of any number less than 1 is negative). LL is calculated through iteration, using maximum likelihood estimation (MLE).

**Assumption 1:** Your dependent variable should be measured on a **dichotomous scale**.

**Assumption 2:** You have one or more independent variables, which can be either continuous (i.e., an interval or ratio variable) or categorical (i.e., an ordinal or nominal variable).

**Assumption 3:** You should have independence of observations and the dependent variable should have mutually exclusive and exhaustive categories.

**Assumption 4:** There needs to be a linear relationship between any continuous independent variables and the logit transformation of the dependent variable.