# Maximum Likelihood Estimation: Numerical Solution for Bernoulli Distribution[1]

The Bernoulli distribution has the probability mass function:

$$f(y|\theta) = \theta^y (1-\theta)^{1-y} \tag{1}$$

where: $y = \{0, 1\}$ and $\theta \in [0, 1]$.

## Maximum Likelihood Estimation

Let $y_1, \ldots, y_n$, denote the data. Assume, $\forall i : y_i$ independente random variables, share the same parameter from a Bernoulli distribution described in (1).

$$f(y_i|\theta) = \theta^{y_i}(1-\theta)^{1-y_i} \tag{2}$$

then, their join probability distribution is:

$$f(y_1, \ldots, y_n|\theta) = \prod_{i=1}^{n} f(y_i|\theta)$$

$$= \prod_{i=1}^{n} \theta^{y_i}(1-\theta)^{1-y_i} \tag{3}$$

The likelihood function is:

$$\mathscr{L}(\theta|y_1, \ldots, y_n) = f(y_1, \ldots, y_n|\theta) \tag{4}$$

The log-likelihood function is:

$$\ell(\theta|y_1, \ldots, y_n) = \log(\mathscr{L}(\theta|y_1, \ldots, y_n))$$

$$= \log \prod_{i=1}^{n} \theta^{y_i}(1-\theta)^{1-y_i}$$

$$= \sum_{i=1}^{n} \log\left(\theta^{y_i}(1-\theta)^{1-y_i}\right)$$

$$= \sum_{i=1}^{n} \left(y_i \log \theta + (1-y_i) \log(1-\theta)\right) \tag{5}$$

The maximum likelihood estimator (MLE), denoted by $\hat{\theta}$, is such that:

$$\hat{\theta} = \operatorname{argmax}\left\{\ell(\theta|y_1, \ldots, y_n)\right\} \tag{6}$$

To maximize the log-likelihood function (5), requires the derivative with respect to $\theta$. The resulted function is called the Score function, denoted by $U(\theta|y_1, \ldots, y_n)$.

$$
\begin{aligned}
U(\theta|y_1, \ldots, y_n) &= \frac{\partial \ell(\theta|y_1, \ldots, y_n)}{\partial \theta} \\
&= \sum_{i=1}^{n} \left( y_i \frac{1}{\theta} + (1 - y_i) \frac{1}{1 - \theta} (-1) \right) \\
&= \frac{1}{\theta(1 - \theta)} \sum_{i=1}^{n} y_i - \frac{n}{1 - \theta}
\end{aligned} \tag{7}
$$

Then, the MLE $\hat{\theta}$ is the solution of:

$$
U\left(\theta = \hat{\theta}|y_1, \ldots, y_n, \lambda\right) = 0 \tag{8}
$$

## Maximum Likelihood Estimation: Newton-Raphson Method

Just for notation, let write equation (8) as:

$$
U(\theta^*) = 0 \tag{9}
$$

The equation (9), generally, is a nonlinear equation, that can be aproximate by Taylor Series:

$$
U(\theta^*) \approx U\left(\theta^{(t)}\right) + U'\left(\theta^{(t)}\right)\left(\theta^* - \theta^{(t)}\right) \tag{10}
$$

Then, using (10) into (9), and solving for $\theta^*$:

$$
\begin{aligned}
U\left(\theta^{(t)}\right) + U'\left(\theta^{(t)}\right)\left(\theta^* - \theta^{(t)}\right) &= 0 \\
\theta^* &= \theta^{(t)} - \frac{U\left(\theta^{(t)}\right)}{U'\left(\theta^{(t)}\right)}
\end{aligned} \tag{11}
$$

where $U'$ is the derivative of the Score function (7) respect of $\theta$.

$$
\begin{aligned}
U'(\theta|y_1, \ldots, y_n, \lambda) &= \frac{\partial U(\theta|y_1, \ldots, y_n, \lambda)}{\partial \theta} \\
&= \frac{2\theta - 1}{\theta^2(1 - \theta)^2} \sum_{i=1}^{n} y_i - \frac{n}{(1 - \theta)^2}
\end{aligned} \tag{12}
$$

Then, with the Newton-Raphson method: starting with an initial guess $\theta^{(1)}$ successive approximations are obtained using (13), until the iterative process converges.

$$
\theta^{(t+1)} = \theta^{(t)} - \frac{U\left(\theta^{(t)}\right)}{U'\left(\theta^{(t)}\right)} \tag{13}
$$

In order to example the use of Newton-Rapshon method, we use the data of total July rainfall (in millimeters) at Quilpie, Australia stored into the data set "quilpie", where the variable y is a dicotomic variable[2].

---

[2]This data was taken from package: GLMsData.

```
# load data
data("quilpie")
Y <- quilpie$y
```

We load the code developed into our R function `MLE_NR_Bernoulli`, stored in the R object with the same name.

```
# load the function to solve by Newton-Raphson
load("MLE_NR_Bernoulli.RData")
```

The function `MLE_NR_Bernoulli` takes $\theta = 0.5$ as a first guess for the iterative process and, besides some other default parameters that can be modified, only needs the data vector Y.

```
# MLE by Newton-Raphson (NR) for Bernoulli distribution
MLE_NR_Bernoulli(Y)
```

```
##       ML Estimator Likelihood              Log-Likelihood
## [1,] "0.5000000"  "3.3881317890172e-21"  "-47.1340082780763"
## [2,] "0.5147059"  "3.48927745358427e-21" "-47.1045922714519"
## [3,] "0.5147059"  "3.48927745358427e-21" "-47.1045922714519"
```

Then, the MLE by Newton-Raphson method: $\hat{\theta} = 0.5147059$.

### Maximum Likelihood Estimation: Fisher-Scoring Method

A distribution belongs to the exponential family if it can be written in the form:

$$f(y|\theta) = \exp\left\{\frac{a(y)\,b(\theta) - c(\theta)}{\phi} + d(y, \phi)\right\} \tag{14}$$

Since (1) can be written as a member of exponential family as in (14):

$$
\begin{aligned}
f(y|\theta, \lambda) &= \exp\left\{\log\left(\theta^y(1-\theta)^{1-y}\right)\right\} \\
&= \exp\left\{y\,\log\left(\frac{\theta}{1-\theta}\right) - (-\log\,(1-\theta))\right\}
\end{aligned} \tag{15}
$$

where, $a(y) = y$, $b(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$, $c(\theta) = -\log\,(1-\theta)$, $\phi = 1$, and $d(y, \phi) = 0$.

Then, since the Bernoulli distribution belongs to the exponential family, it can be show that the variance of $U$, denoted by $\mathcal{J}$, is:

$$\mathcal{J} = \text{Var}\{U\} = -\text{E}\{U'\} \tag{16}$$

where:

$$\text{E}\{U'\} = -\frac{1}{\phi}\left(b''(\theta)\frac{c'(\theta)}{b'(\theta)} - c''(\theta)\right) \tag{17}$$

For MLE, it is common to approximate $U'$ by its expected value $\mathsf{E}\{U'\}$. In this case:

$$
\begin{aligned}
\mathcal{J} &= -\mathsf{E}\{U'\} \\
&= \mathsf{E}\{-U'\} \\
&= \mathsf{E}\left\{-\sum_{i=1}^{n} U'_i\right\} \\
&= \sum_{i=1}^{n} -\mathsf{E}\{U'_i\} \\
&= \sum_{i=1}^{n} -\frac{1}{\phi}\left(b''(\theta)\frac{c'(\theta)}{b'(\theta)} - c''(\theta)\right)
\end{aligned}
\tag{18}
$$

where, using (1), the previous derivaties:

$$
\begin{aligned}
b'(\theta) &= \frac{1}{\theta(1-\theta)} \\
b''(\theta) &= \frac{2\theta-1}{\theta^2(1-\theta)^2} \\
c'(\theta) &= \frac{1}{1-\theta} \\
c''(\theta) &= \frac{1}{(1-\theta)^2} \\
\frac{1}{\phi} &= 1
\end{aligned}
$$

Then, replacing them into (18):

$$
\begin{aligned}
\mathcal{J} &= \sum_{i=1}^{n} -\left\{\frac{2\theta-1}{\theta^2(1-\theta)^2}\frac{\frac{1}{1-\theta}}{\frac{1}{\theta(1-\theta)}} - \frac{1}{(1-\theta)^2}\right\} \\
&= \sum_{i=1}^{n} \frac{1}{\theta(1-\theta)} \\
&= \frac{n}{\theta(1-\theta)}
\end{aligned}
\tag{19}
$$

Finally:

$$
\begin{aligned}
\mathcal{J} &= -\mathsf{E}\{U'\} = \frac{n}{\theta(1-\theta)} \\
-\mathcal{J} &= \mathsf{E}\{U'\} = \frac{n}{\theta(1-\theta)}
\end{aligned}
\tag{20}
$$

Then, approximating $U'$ by its expected value $\mathsf{E}\{U'\}$, the equation (13) results into:

$$
\theta^{(t+1)} = \theta^{(t)} + \frac{U\left(\theta^{(t)}\right)}{\mathcal{J}\left(\theta^{(t)}\right)}
\tag{21}
$$

In order to example the use of Fisher-Scoring method, we use the same data used in the Newton-Rapshon method. We load the code developed into our R function `MLE_FS_Bernoulli`, stored in the R object with the same name.

```r
# load the function to solve by Fisher-Scoring
load("MLE_FS_Bernoulli.RData")
```

The function `MLE_FS_Bernoulli` takes $\theta = 0.5$ as a first guess for the iterative process and, besides some other default parameters that can be modified, only needs the data vector Y.

```r
# MLE by Fisher-Scoring (FS) for Bernoulli distribution
MLE_FS_Bernoulli(Y)

##       ML Estimator Likelihood           Log-Likelihood
## [1,] "0.5000000"  "3.3881317890172e-21" "-47.1340082780763"
## [2,] "0.5147059"  "3.48927745358427e-21" "-47.1045922714519"
## [3,] "0.5147059"  "3.48927745358427e-21" "-47.1045922714519"
```

Then, the MLE by Fisher-Scoring method: $\hat{\theta} = 0.5147059$.

In summary, the same estimate for the MLE is achieved by both approaches: the Newton-Raphson and the Fisher-Scoring method.

### Naive Approach

The main idea behind the Maximum Likelihood (ML) method is to choose those estimates for the unknown parameters that maximize the join probability of our observed data (our sample). Keeping in mind this idea, if we want to get the MLE and avoiding to implement a numerical solution, a naive approach is to set a large range of possible values for unknown parameters, evaluate the log-likelihood function (also, the likelihood function) at each point and the point for which the log-likelihood function (also, the likelihood function) reaches its maximum value, will be our MLE looking for.
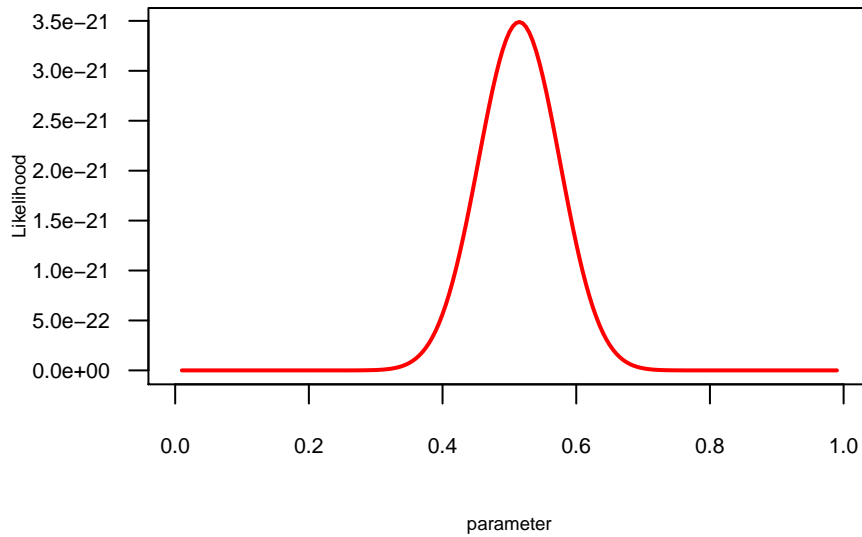
We set a set of values por the parameter, from $\hat{\theta} = 0.01$ to $\hat{\theta} = 0.99$, spaced by 0.005; and evaluate (4) and (5) at each point from the set of values.

```r
# set a large range of values for the parameter
theta <- seq(0.01, 0.99, 0.005)
tabla <- matrix(c(NA,NA,NA), nrow = length(theta), ncol = 3)
tabla[,1] <- theta

# evaluate the likelihood and the log-likelihood at each point
for (i in 1:length(theta)){
  tabla[i,2] <- prod(dbinom(Y, size = 1, prob = theta[i], log = FALSE))
  tabla[i,3] <- sum(dbinom(Y, size = 1, prob = theta[i], log = TRUE))
}
colnames(tabla) <- c("theta", "Likelihood", "Log-Likelihood")
df_tabla <- as.data.frame(tabla)
```

Then, the plot of likelihood function evaluated at each point:

**Figure 1: Likelihooh Function**
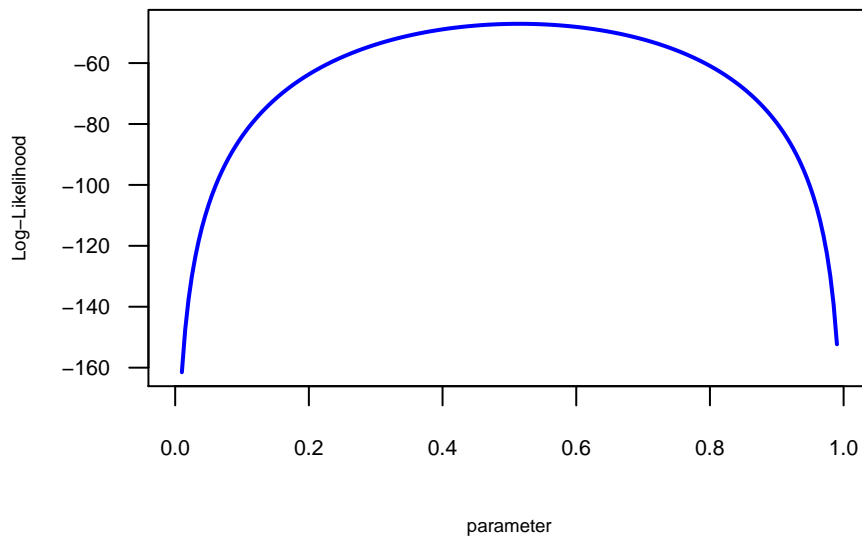


The point for which the likelihood function is maximum, is our MLE from a naive approach:

```
mxl <- max(tabla[,2])
df_tabla %>% filter(`Likelihood` == mxl)

##    theta   Likelihood Log-Likelihood
## 1 0.515 3.489236e-21       -47.1046
```

Also, the plot of log-likelihood function evaluated at each point:

**Figure 2: Log–Likelihooh Function**

The point for which the log-likelihood function is maximum, which is the same point at the likelihood function reaches its maximum value, is our MLE from a naive approach:

```
mxllog <- max(tabla[,3])
df_tabla %>% filter(`Log-Likelihood` == mxllog)

##    theta    Likelihood Log-Likelihood
## 1 0.515 3.489236e-21       -47.1046
```

As we can see, using this naive approach, we reach a value that is close enough to that which is reached using the numerical solution: the Newton-Raphson and the Fisher-Scoring method.