

1.

$$f(y_i; \beta) = \mu_i^{y_i} (1 - \mu_i)^{1-y_i}, \quad y_i \in \{0, 1\}, \quad \text{where } \mu_i = \Phi(\mathbf{x}_i^\top \beta).$$

So each term can be written as

$$P(y_i | x_i, \beta) = (\Phi(x_i^\top \beta))^{y_i} (1 - \Phi(x_i^\top \beta))^{1-y_i}$$

The likelihood function is

$$L(\beta) = \prod_{i=1}^n P(y_i | x_i, \beta), \quad y_i \in \{0, 1\}$$

Let

$$P(y_i = 1 | x_i, \beta) = \Phi(x_i^\top \beta), \quad P(y_i = 0 | x_i, \beta) = 1 - \Phi(x_i^\top \beta), \quad y_i \in \{0, 1\}.$$

The Product likelihood:

$$L(\beta) = \prod_{i=1}^n P(y_i | x_i, \beta)$$

$$\Rightarrow L(\beta) = \prod_{i=1}^n (\Phi(x_i^\top \beta))^{y_i} (1 - \Phi(x_i^\top \beta))^{1-y_i}$$

$$\Rightarrow \log L(\beta) = \log \prod_{i=1}^n \left((\Phi(x_i^\top \beta))^{y_i} (1 - \Phi(x_i^\top \beta))^{1-y_i} \right) \quad (\text{taking log})$$

$$\Rightarrow \log L(\beta) = \sum_{i=1}^n \log \left((\Phi(x_i^\top \beta))^{y_i} (1 - \Phi(x_i^\top \beta))^{1-y_i} \right) \quad (\text{since } \log \prod_i a_i = \sum_i \log a_i)$$

$$\Rightarrow \log L(\beta) = \sum_{i=1}^n [\log (\Phi(x_i^\top \beta))^{y_i} + \log ((1 - \Phi(x_i^\top \beta))^{1-y_i})] \quad (\text{since } \log(ab) = \log(a) + \log(b))$$

$$\Rightarrow \log L(\beta) = \sum_{i=1}^n [y_i \log \Phi(x_i^\top \beta) + (1 - y_i) \log (1 - \Phi(x_i^\top \beta))] \quad (\text{since } \log(a^b) = b \log(a)).$$

Sanity check using $y_i \in \{0, 1\}$:

- If $y_i = 1$, the term becomes $\log \Phi(\mathbf{x}_i^\top \beta)$
- If $y_i = 0$, the term becomes $\log (1 - \Phi(\mathbf{x}_i^\top \beta))$

The log-likelihood function for probit regression represents the logarithm of the probability of observing the entire dataset given parameters β . We maximize this to find the best-fitting model:

$$\log L(\beta) = \sum_{i=1}^n [y_i \log \Phi(\mathbf{x}_i^\top \beta) + (1 - y_i) \log (1 - \Phi(\mathbf{x}_i^\top \beta))]$$

The predicted probability $\mu(\mathbf{x})$ is obtained by transforming the linear predictor $\beta^\top \mathbf{x}$ through the standard normal CDF Φ . This maps any real value to a probability between 0 and 1:

$$\mu(\mathbf{x}) = \Phi(\beta^\top \mathbf{x}) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\beta^\top \mathbf{x}} \exp\left(-\frac{u^2}{2}\right) du$$

The derivative of the standard normal CDF with respect to its argument is the standard normal PDF $\phi(z)$.

$$\frac{d}{dz} \Phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) = \phi(z)$$

Applying this to the predicted probability function, the rate of change of $\mu(\mathbf{x})$ with respect to the linear predictor is:

$$\frac{\partial}{\partial(\beta^\top \mathbf{x})} \mu(\mathbf{x}) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x})^2}{2}\right)$$

Let $g_i(\beta)$ be the contribution of the i -th observation to the log-likelihood. To compute the gradient efficiently, we first consider the log-likelihood contribution from a single observation i . This allows to derive the gradient term-by-term:

$$g_i(\beta) = y_i \log[\Phi(\beta^\top \mathbf{x}_i)] + (1 - y_i) \log[1 - \Phi(\beta^\top \mathbf{x}_i)]$$

We compute the gradient of the term that activates when the outcome is positive ($y_i = 1$). This uses the chain rule three times: for the logarithm, the CDF, and the linear predictor:

$$\begin{aligned} \frac{\partial}{\partial \beta} \left\{ y_i \log[\Phi(\beta^\top \mathbf{x}_i)] \right\} &= y_i \cdot \frac{1}{\Phi(\beta^\top \mathbf{x}_i)} \cdot \frac{\partial}{\partial \beta} \Phi(\beta^\top \mathbf{x}_i) \\ &\quad \text{(Step 1: Apply chain rule; derivative of } \log(u) \text{ is } \frac{1}{u}) \\ &= y_i \cdot \frac{1}{\Phi(\beta^\top \mathbf{x}_i)} \cdot \phi(\beta^\top \mathbf{x}_i) \cdot \mathbf{x}_i \\ &\quad \text{(Step 2: Use } \frac{d\Phi(z)}{dz} = \phi(z) \text{ and } \frac{\partial(\beta^\top \mathbf{x}_i)}{\partial \beta} = \mathbf{x}_i) \\ &= y_i \cdot \frac{1}{\Phi(\beta^\top \mathbf{x}_i)} \cdot \left(\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \right) \cdot \mathbf{x}_i \\ &\quad \text{(Step 3: Substitute the explicit form of the standard normal PDF)} \end{aligned}$$

Similarly, we compute the gradient of the term that activates when the outcome is negative ($y_i = 0$). The key difference is the negative sign from differentiating $(1 - \Phi)$:

$$\begin{aligned} \frac{\partial}{\partial \beta} \left\{ (1 - y_i) \log[1 - \Phi(\beta^\top \mathbf{x}_i)] \right\} &= (1 - y_i) \cdot \frac{1}{1 - \Phi(\beta^\top \mathbf{x}_i)} \cdot \frac{\partial}{\partial \beta} (1 - \Phi(\beta^\top \mathbf{x}_i)) \\ &\quad \text{(Step 1: Apply chain rule to } \log(1 - u) \text{)} \\ &= (1 - y_i) \cdot \frac{1}{1 - \Phi(\beta^\top \mathbf{x}_i)} \cdot (-\phi(\beta^\top \mathbf{x}_i)) \cdot \mathbf{x}_i \\ &\quad \text{(Step 2: Derivative of } (1 - \Phi) \text{ is } -\phi; \text{ include gradient of linear predictor)} \\ &= (1 - y_i) \cdot \frac{1}{1 - \Phi(\beta^\top \mathbf{x}_i)} \cdot \left(-\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \right) \cdot \mathbf{x}_i \\ &\quad \text{(Step 3: Substitute the explicit form of the standard normal PDF)} \\ &= -(1 - y_i) \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right)}{1 - \Phi(\beta^\top \mathbf{x}_i)} \cdot \mathbf{x}_i \\ &\quad \text{(Step 4: Rearrange to show the PDF in the numerator)} \end{aligned}$$

Combining both partial derivatives gives the complete gradient contribution from observation i :

$$\frac{\partial g_i}{\partial \beta} = \left(y_i \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right)}{\Phi(\beta^\top \mathbf{x}_i)} - (1 - y_i) \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right)}{1 - \Phi(\beta^\top \mathbf{x}_i)} \right) \mathbf{x}_i$$

Starting with the gradient contributions from all n observations, we sum the individual derivatives we computed earlier:

$$\nabla_{\beta} \log L(\beta) = \sum_{i=1}^n \left[y_i \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right)}{\Phi(\beta^\top \mathbf{x}_i)} - (1 - y_i) \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right)}{1 - \Phi(\beta^\top \mathbf{x}_i)} \right] \mathbf{x}_i$$

(This is the sum of all $\frac{\partial g_i}{\partial \beta}$ terms)

Both terms in the brackets contain the same normal PDF $\phi(\beta^\top \mathbf{x}_i) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right)$. We factor this out:

$$= \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \left[\frac{y_i}{\Phi(\beta^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\beta^\top \mathbf{x}_i)} \right] \mathbf{x}_i$$

(Factoring out the standard normal PDF $\phi(\beta^\top \mathbf{x}_i)$)

To simplify the expression in brackets, we find a common denominator and combine the fractions:

$$= \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \cdot \frac{y_i(1 - \Phi(\beta^\top \mathbf{x}_i)) - (1 - y_i)\Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \cdot \mathbf{x}_i$$

(Common denominator is $\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))$)

Expanding the numerator:

$$\begin{aligned} \text{Numerator} &= y_i(1 - \Phi(\beta^\top \mathbf{x}_i)) - (1 - y_i)\Phi(\beta^\top \mathbf{x}_i) \\ &= y_i - y_i\Phi(\beta^\top \mathbf{x}_i) - \Phi(\beta^\top \mathbf{x}_i) + y_i\Phi(\beta^\top \mathbf{x}_i) \\ &= y_i - \Phi(\beta^\top \mathbf{x}_i) \end{aligned}$$

(The $y_i\Phi(\beta^\top \mathbf{x}_i)$ terms cancel)

Therefore:

$$= \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \cdot \frac{y_i - \Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \cdot \mathbf{x}_i$$

The gradient of the log-likelihood is:

$$\nabla_{\beta} \log L(\beta) = \sum_{i=1}^n \phi(\beta^\top \mathbf{x}_i) \cdot \frac{y_i - \Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \cdot \mathbf{x}_i$$

where $\phi(\cdot)$ is the standard normal PDF and $\Phi(\cdot)$ is the standard normal CDF.

To find the Hessian matrix, we take the derivative of the gradient with respect to β . Since each term in the gradient has the form [scalar] $\cdot \mathbf{x}_i$, the Hessian will involve outer products $\mathbf{x}_i \mathbf{x}_i^\top$:

$$\nabla_{\beta}^2 \log L(\beta) = \sum_{i=1}^n \left(\frac{\partial}{\partial \beta} \left[\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \left(\frac{y_i}{\Phi(\beta^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\beta^\top \mathbf{x}_i)} \right) \right] \right) \mathbf{x}_i^\top$$

We need to compute the derivative of the product of the PDF term and the fraction term.

First, we compute how the normal PDF changes with respect to its argument. Using the chain rule on the exponential:

$$\begin{aligned} \frac{d}{d(\beta^\top \mathbf{x}_i)} \left[\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \right] &= -(\beta^\top \mathbf{x}_i) \cdot \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2}\right) \\ &= -(\beta^\top \mathbf{x}_i) \cdot \phi(\beta^\top \mathbf{x}_i) \end{aligned}$$

Using the quotient rule $\frac{d}{dz} \left(\frac{u}{v} \right) = -\frac{u \cdot v'}{v^2}$ when u is constant and v depends on z :

$$\begin{aligned} \frac{d}{d(\boldsymbol{\beta}^\top \mathbf{x}_i)} \left(\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} \right) &= -y_i \cdot \frac{\phi(\boldsymbol{\beta}^\top \mathbf{x}_i)}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} \\ &= -y_i \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}\right)}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} \end{aligned}$$

Similarly, for the second term, but now the derivative of the denominator $(1 - \Phi)$ is $-\phi$, giving a positive result:

$$\begin{aligned} \frac{d}{d(\boldsymbol{\beta}^\top \mathbf{x}_i)} \left(\frac{1 - y_i}{1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} \right) &= (1 - y_i) \cdot \frac{\phi(\boldsymbol{\beta}^\top \mathbf{x}_i)}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \\ &= (1 - y_i) \cdot \frac{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}\right)}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \end{aligned}$$

Adding the derivatives of both fraction terms and factoring out the common PDF:

$$\begin{aligned} \frac{d}{d(\boldsymbol{\beta}^\top \mathbf{x}_i)} \left(\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} \right) &= -y_i \cdot \frac{\phi(\boldsymbol{\beta}^\top \mathbf{x}_i)}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} + (1 - y_i) \cdot \frac{\phi(\boldsymbol{\beta}^\top \mathbf{x}_i)}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \\ &\quad \text{(Substituting the derivatives computed above)} \\ &= \phi(\boldsymbol{\beta}^\top \mathbf{x}_i) \left[-\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} + \frac{1 - y_i}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \right] \\ &\quad \text{(Factoring out the common normal PDF)} \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}\right) \left[-\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} + \frac{1 - y_i}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \right] \\ &\quad \text{(Writing out the PDF explicitly)} \end{aligned}$$

To find the second derivative, we apply the product rule to the product of the PDF and the fraction terms. Recall that for a product $f \cdot g$, we have $(f \cdot g)' = f' \cdot g + f \cdot g'$:

$$\begin{aligned} \frac{d}{d(\boldsymbol{\beta}^\top \mathbf{x}_i)} \left[\phi(\boldsymbol{\beta}^\top \mathbf{x}_i) \cdot \left(\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} \right) \right] & \\ \text{(Setting up the product rule for [PDF] } \times \text{ [fraction terms])} & \end{aligned}$$

The product rule gives two terms: the derivative of the PDF times the original fractions, plus the PDF times the derivative of the fractions:

$$\begin{aligned} &= \left(-(\boldsymbol{\beta}^\top \mathbf{x}_i) \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}} \right) \left(\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} \right) \\ &\quad + \frac{1}{\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}} \left[-\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} + \frac{1 - y_i}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \right] \\ &\quad \text{(First term: derivative of } \phi, \text{ which is } \phi'(z) = -z\phi(z)) \\ &\quad \text{(Second term: } \phi \text{ times the derivative of fractions from the previous calculation)} \end{aligned}$$

We factor out $\phi(\boldsymbol{\beta}^\top \mathbf{x}_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}}$ from both terms:

$$\begin{aligned} &= \frac{1}{\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}} \left[-(\boldsymbol{\beta}^\top \mathbf{x}_i) \left(\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)} \right) \right. \\ &\quad \left. + \frac{1}{\sqrt{2\pi}} e^{-\frac{(\boldsymbol{\beta}^\top \mathbf{x}_i)^2}{2}} \left(-\frac{y_i}{\Phi(\boldsymbol{\beta}^\top \mathbf{x}_i)^2} + \frac{1 - y_i}{(1 - \Phi(\boldsymbol{\beta}^\top \mathbf{x}_i))^2} \right) \right] \\ &\quad \text{(Factoring out the first PDF; note the second term has } \phi^2) \end{aligned}$$

We convert the fractions to common denominators. For the first bracket, we use the fact that:

$$\frac{y_i}{\Phi} - \frac{1 - y_i}{1 - \Phi} = \frac{y_i(1 - \Phi) - (1 - y_i)\Phi}{\Phi(1 - \Phi)} = \frac{y_i - \Phi}{\Phi(1 - \Phi)}$$

For the second bracket, we find a common denominator of $\Phi^2(1 - \Phi)^2$:

$$= \phi(\beta^\top \mathbf{x}_i) \left[-(\beta^\top \mathbf{x}_i) \cdot \frac{y_i - \Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \right. \\ \left. - \phi(\beta^\top \mathbf{x}_i) \cdot \frac{y_i(1 - \Phi(\beta^\top \mathbf{x}_i))^2 - (1 - y_i)\Phi(\beta^\top \mathbf{x}_i)^2}{\Phi(\beta^\top \mathbf{x}_i)^2(1 - \Phi(\beta^\top \mathbf{x}_i))^2} \right]$$

(Simplifying fractions; note the negative sign in the second term)

The Hessian matrix is the matrix of second partial derivatives. We start by differentiating the gradient with respect to β :

$$\nabla_{\beta}^2 \log L(\beta) = \sum_{i=1}^n \left(\frac{\partial}{\partial \beta} \left[\frac{1}{\sqrt{2\pi}} \exp \left(-\frac{(\beta^\top \mathbf{x}_i)^2}{2} \right) \left(\frac{y_i}{\Phi(\beta^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\beta^\top \mathbf{x}_i)} \right) \right] \right) \mathbf{x}_i^\top$$

(The outer product \mathbf{x}_i^\top comes from differentiating the \mathbf{x}_i factor in the gradient)

Using the product rule calculation from the previous work, we substitute the derivative and include the outer product $\mathbf{x}_i \mathbf{x}_i^\top$:

$$= \sum_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(\beta^\top \mathbf{x}_i)^2}{2}} \left[-(\beta^\top \mathbf{x}_i) \left(\frac{y_i}{\Phi(\beta^\top \mathbf{x}_i)} - \frac{1 - y_i}{1 - \Phi(\beta^\top \mathbf{x}_i)} \right) \right. \\ \left. + \frac{1}{\sqrt{2\pi}} e^{-\frac{(\beta^\top \mathbf{x}_i)^2}{2}} \left(-\frac{y_i}{\Phi(\beta^\top \mathbf{x}_i)^2} + \frac{1 - y_i}{(1 - \Phi(\beta^\top \mathbf{x}_i))^2} \right) \right] \mathbf{x}_i \mathbf{x}_i^\top$$

(Note: $\mathbf{x}_i \mathbf{x}_i^\top$ is a $p \times p$ matrix where p is the dimension of β)

We factor out the negative sign and rearrange terms to highlight the structure. Note that we've corrected the sign in the numerator:

$$= - \sum_{i=1}^n \phi(\beta^\top \mathbf{x}_i) \left[\phi(\beta^\top \mathbf{x}_i) \cdot \frac{y_i(1 - \Phi(\beta^\top \mathbf{x}_i))^2 + (1 - y_i)\Phi(\beta^\top \mathbf{x}_i)^2}{\Phi(\beta^\top \mathbf{x}_i)^2(1 - \Phi(\beta^\top \mathbf{x}_i))^2} \right. \\ \left. + (\beta^\top \mathbf{x}_i) \cdot \frac{y_i - \Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \right] \mathbf{x}_i \mathbf{x}_i^\top$$

(Factoring out $-\phi$ and using simplified fractions from earlier)

Distributing the negative sign to both terms inside the brackets gives the final expression:

$$= \sum_{i=1}^n \phi(\beta^\top \mathbf{x}_i) \left[-\phi(\beta^\top \mathbf{x}_i) \cdot \frac{y_i(1 - \Phi(\beta^\top \mathbf{x}_i))^2 + (1 - y_i)\Phi(\beta^\top \mathbf{x}_i)^2}{\Phi(\beta^\top \mathbf{x}_i)^2(1 - \Phi(\beta^\top \mathbf{x}_i))^2} \right. \\ \left. - (\beta^\top \mathbf{x}_i) \cdot \frac{y_i - \Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \right] \mathbf{x}_i \mathbf{x}_i^\top$$

$$\nabla_{\beta}^2 \log L(\beta) = \sum_{i=1}^n \phi(\beta^\top \mathbf{x}_i) \left[-\phi(\beta^\top \mathbf{x}_i) \cdot \frac{y_i(1 - \Phi(\beta^\top \mathbf{x}_i))^2 + (1 - y_i)\Phi(\beta^\top \mathbf{x}_i)^2}{\Phi(\beta^\top \mathbf{x}_i)^2(1 - \Phi(\beta^\top \mathbf{x}_i))^2} - (\beta^\top \mathbf{x}_i) \cdot \frac{y_i - \Phi(\beta^\top \mathbf{x}_i)}{\Phi(\beta^\top \mathbf{x}_i)(1 - \Phi(\beta^\top \mathbf{x}_i))} \right] \mathbf{x}_i \mathbf{x}_i^\top$$

```

1 def newton_raphson(model, tol=1e-5, max_iter=1000, display=True):
2     i = 0
3     error = 100
4     if display:
5         header = f'{"Iteration_k":<13>{"Log-likelihood":<16>{"theta":<60}'

```

```

6     print(header)
7     print("-" * len(header))
8
9     while np.any(error > tol) and i < max_iter:
10         H, G = model.H(), model.G()
11          $\beta_{\text{new}} = \text{model}.\beta - (\text{np.linalg.inv}(H) @ G)$ 
12         error =  $\beta_{\text{new}} - \text{model}.\beta$ 
13         model. $\beta = \beta_{\text{new}}$ 
14         if display:
15              $\beta_{\text{list}} = [f'\{t:.3\}' \text{ for } t \text{ in list}(\text{model}.\beta.\text{flatten}())]$ 
16             update = f'\{i:<13\}{model.logL():<16.8}\{ $\beta_{\text{list}}$ \}'
17             print(update)
18
19         i += 1
20
21     print(f'Number of iterations: {i}')
22     print(f' $\beta_{\text{hat}} = \{\text{model}.\beta.\text{flatten}()\}$ ')
23     return model. $\beta.\text{flatten}()$ 

```

```

1 class ProbitModel:
2     def __init__(self, X, y, beta_init):
3         self.X = np.asarray(X, dtype=float)
4         self.y = np.asarray(y, dtype=float).reshape(-1)
5         self. $\beta$  = np.asarray(beta_init, dtype=float).reshape(-1)
6
7     def logL(self) -> float:
8         z = self.X @ self. $\beta$ 
9         return float(np.sum(self.y * norm.logcdf(z) + (1 - self.y) * norm.logsf(z)))
10
11     def G(self) -> np.ndarray:
12         z = self.X @ self. $\beta$ 
13         Phi = norm.cdf(z)
14         phi = norm.pdf(z)
15         eps = np.finfo(float).eps
16         denom = np.clip(Phi * (1 - Phi), eps, None)
17         w = phi * (self.y - Phi) / denom
18         return self.X.T @ w
19
20     def H(self) -> np.ndarray:
21         z = self.X @ self. $\beta$ 
22         Phi = norm.cdf(z)
23         phi = norm.pdf(z)
24         eps = np.finfo(float).eps
25         PhiS = np.clip(Phi, eps, 1 - eps)
26         OMS = 1 - PhiS
27         A_num = self.y * (OMS**2) + (1 - self.y) * (PhiS**2)
28         A_den = (PhiS**2) * (OMS**2)
29         A = A_num / A_den
30         B = (self.y - PhiS) / (PhiS * OMS)
31         coeff = phi * ( - phi * A - z * B )
32         return self.X.T @ (coeff[:, None] * self.X)

```

2.

$$\mathbf{X} = \begin{bmatrix} 1 & 2 & 4 \\ 1 & 1 & 1 \\ 1 & 4 & 3 \\ 1 & 5 & 6 \\ 1 & 3 & 5 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \quad \beta_{(0)} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$

Based on the given dataset and the initial value of $\beta_{(0)} = (0.1, 0.1, 0.1)^\top$, we applied the Newton–Raphson algorithm to maximize the probit log-likelihood function. The iteration process converged after 5 steps, yielding the following maximum likelihood estimate:

$$\hat{\beta} = \begin{bmatrix} -1.5463 \\ 0.7778 \\ -0.0971 \end{bmatrix}$$

```
Iteration_k  Log-likelihood  θ
-----
0            -2.3796884    ['-1.34', '0.775', '-0.157']
1            -2.3687526    ['-1.53', '0.775', '-0.0981']
2            -2.3687294    ['-1.55', '0.778', '-0.0971']
3            -2.3687294    ['-1.55', '0.778', '-0.0971']
4            -2.3687294    ['-1.55', '0.778', '-0.0971']
Number of iterations: 5
β_hat = [-1.54625858  0.77778952 -0.09709757]
```

Optimization terminated successfully.
Current function value: 0.473746
Iterations 6

Probit Regression Results

```
=====
Dep. Variable:          y      No. Observations:      5
Model:                  Probit  Df Residuals:          2
Method:                  MLE    Df Model:          2
Date:                   Mon, 29 Sep 2025  Pseudo R-squ.:    0.2961
Time:                   19:55:22  Log-Likelihood:    -2.3687
converged:              True     LL-Null:          -3.3651
Covariance Type:        nonrobust  LLR p-value:       0.3692
=====
```

	coef	std err	z	P> z	[0.025	0.975]
const	-1.5463	1.866	-0.829	0.407	-5.204	2.111
x1	0.7778	0.788	0.986	0.324	-0.768	2.323
x2	-0.0971	0.590	-0.165	0.869	-1.254	1.060

```
=====
```

Estimated coefficients (β_{hat}):
[-1.54625858 0.77778952 -0.09709757]

```
1 X = np.array([
2     [1, 2, 4],
3     [1, 1, 1],
4     [1, 4, 3],
5     [1, 5, 6],
6     [1, 3, 5],
7 ], dtype=float)
8
9 y = np.array([1, 0, 1, 1, 0], dtype=float)
10 beta0 = np.array([0.1, 0.1, 0.1], dtype=float)
11
12 model = ProbitModel(X, y, beta0)
13 beta_hat = newton_raphson(model, display=True)
```

```
1 import numpy as np
2 import pandas as pd
3 from statsmodels.discrete.discrete_model import Probit
4 import statsmodels.api as sm
5
6 X = np.array([
7     [1, 2, 4],
8     [1, 1, 1],
```



```

9      [1, 4, 3],
10     [1, 5, 6],
11     [1, 3, 5],
12 ], dtype=float)
13
14 y = np.array([1, 0, 1, 1, 0], dtype=float)
15
16 model = Probit(y, X)
17 result = model.fit(displ=True)
18
19 print(result.summary())
20 print("\nEstimated coefficients (beta_hat):")
21 print(result.params)

```

3.

```

Iteration_k  Log-likelihood  θ
-----
0            -5603.7745      ['1.78', '0.74', '0.176']
1            -5239.8158      ['1.58', '0.741', '0.155']
2            -5234.0225      ['1.56', '0.742', '0.148']
3            -5234.0204      ['1.56', '0.742', '0.148']
4            -5234.0204      ['1.56', '0.742', '0.148']
Number of iterations: 5
β_hat = [1.55904334 0.74199832 0.14767199]

```

```

Intercept: 1.5590425134556154
Coefficients: [0.74199947 0.14767126]
Model Score: 0.39440825681225433

```

Optimization terminated successfully.
 Current function value: 5.234020
 Iterations 6

Poisson Regression Results

```
=====
Dep. Variable:          y      No. Observations:      1000
Model:                Poisson  Df Residuals:          997
Method:                MLE     Df Model:            2
Date:                  Mon, 29 Sep 2025  Pseudo R-squ.:      0.3092
Time:                  20:59:09  Log-Likelihood:     -5234.0
converged:              True     LL-Null:           -7577.0
Covariance Type:        nonrobust  LLR p-value:        0.000
=====
```

	coef	std err	z	P> z	[0.025	0.975]
const	1.5590	0.026	60.033	0.000	1.508	1.610
x1	0.7420	0.011	65.799	0.000	0.720	0.764
x2	0.1477	0.010	14.078	0.000	0.127	0.168

The Newton-Raphson algorithm shows good convergence behavior, achieving optimal parameter estimates in only 5 iterations with $\hat{\beta} = [1.55904334, 0.74199832, 0.14767199]^T$. The convergence pattern shows rapid improvement, with the log-likelihood value moving from -5603.7745 at initialization to -5239.8158 after the first iteration, then fine-tuning to -5234.0204 by iteration 3 where practical convergence was achieved. The sklearn implementation produced identical results with intercept 1.55904251 and coefficients $[0.74199947, 0.14767126]$, differing by less than 10^{-6} from Newton-Raphson, while achieving a D^2 model score of 0.3944 indicating that 39.4% of the deviance is explained by the model.

The statsmodels output, converging in 6 iterations (one more than Newton-Raphson's 5) to an identical log-likelihood of -5234.0 , validates the Newton-Raphson implementation with parameter estimates matching to within numerical precision: statsmodels reports $\beta_0 = 1.5590$, $\beta_1 = 0.7420$, and $\beta_2 = 0.1477$, which agree with Newton-Raphson's estimates to 4 decimal places.

The maximum absolute difference between all three methods is less than 1.15×10^{-6} , confirming that Newton-Raphson, sklearn, and statsmodels have converged to the same global optimum. While sklearn reports a D^2 model score of 0.3944 and statsmodels provides a pseudo R -squared of 0.3092 , both indicate strong model fit, and all three implementations achieve the identical maximized log-likelihood of -5234.02 , demonstrating agreement in the optimization solution using different algorithms.

```
1 from scipy.special import gammaln
2
3 class PoissonRegression:
4
5     def __init__(self, y, X, beta):
6         self.X = X
7         self.n, self.k = X.shape
8         self.y = y.reshape(self.n,1)
9         self.beta = beta.reshape(self.k,1)
10
11     def mu(self):
12         return np.exp(self.X @ self.beta)
```

```

14     def logL(self):
15         y = self.y
16          $\mu$  = self. $\mu$ ()
17         return np.sum(y * np.log( $\mu$ ) -  $\mu$  - gammaln(y + 1))
18
19     def G(self):
20         X = self.X
21         y = self.y
22          $\mu$  = self. $\mu$ ()
23         return X.T @ (y -  $\mu$ )
24
25     def H(self):
26         X = self.X
27          $\mu$  = self. $\mu$ ()
28         return -X.T @ ( $\mu$  * X)

```

```

1 X = df[df.columns.drop('y')]
2 y = df['y']
3 Xc=sm.add_constant(X, prepend=True)
4
5 X_m = Xc.values
6 y_m = y.values

```

```

1 init_ $\beta$  = np.array([np.log(y_m.mean()), 0.0, 0.0])
2 poi = PoissonRegression(y_m, X_m, init_ $\beta$ )
3  $\beta$ _hat = newton_raphson(poi)

```

```

1 model = PoissonRegressor(alpha=0, max_iter=1000)
2 model.fit(X, y)
3 print(f"Intercept: {model.intercept_}")
4 print(f"Coefficients: {model.coef_}")
5 print(f"Model Score: {model.score(X, y)}")

```

```

1 Xc=sm.add_constant(X, prepend=True)
2 mod_p = sm.Poisson(y, Xc)
3 res_p = mod_p.fit()
4 print(res_p.summary())

```