

Recall that we specified a glm for binary data as

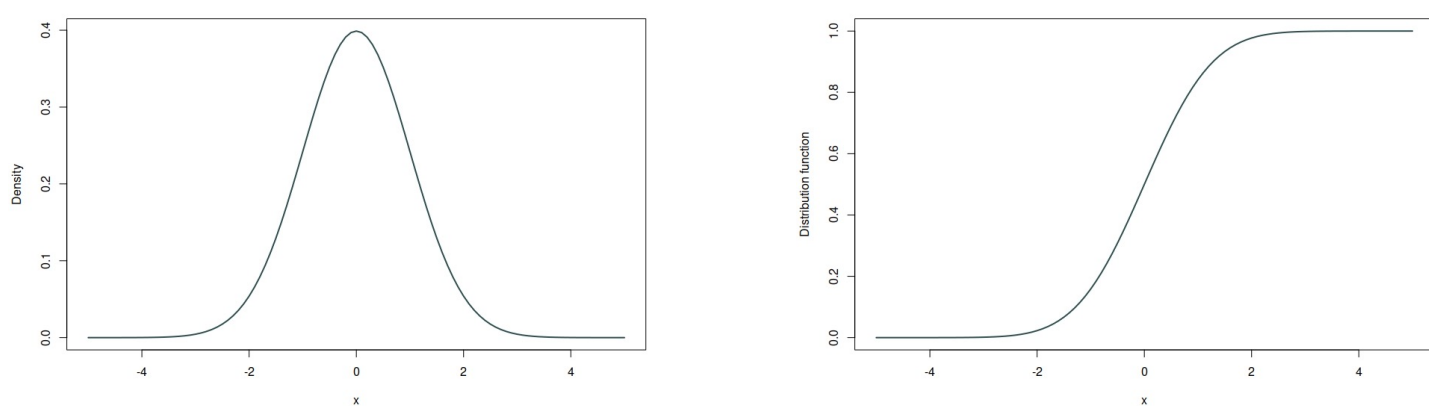
1. $Y_i \sim \text{Bernoulli}(\pi_i)$ independent $i=1, \dots, n$
hence $\pi_i = \mathbb{E}[Y_i] = \mathbb{P}(Y_i=1)$, $\pi_i \in [0,1]$
2. $\eta_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} = \underline{x}_i^T \underline{\beta}$
3. $g(\pi_i) = \eta_i$

We analyzed the case where $g(\cdot)$ is the canonical link function: logit model
However, g could be any function that maps $[0,1] \rightarrow \mathbb{R}$, invertible (and differentiable).
→ cumulative distribution functions are good candidates.

• INTERPRETATION as THRESHOLD MODEL

Assume that $Y_i \sim \text{Bernoulli}(\pi_i)$ $i=1, \dots, n$ and
 $\pi_i = F(\underline{x}_i^T \underline{\beta})$ with F the cdf of a r.v. with distribution SYMMETRIC around zero
Then the regression for Y_i has an interpretation in terms of a model on a CONTINUOUS LATENT (= unobserved) r.v. Y_i^* .

Let us consider, for example, the PROBIT MODEL, where $F = \Phi$ is the cumulative distribution function of a standard Gaussian distribution:



PROBIT REGRESSION: model assumptions

- $Y_i \sim \text{Bernoulli}(\pi_i)$ $i=1, \dots, n$ independent
- $\eta_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} = \underline{x}_i^T \underline{\beta}$
- $g(\pi_i) = \Phi^{-1}(\pi_i) = \eta_i$
⇒ $\pi_i = \Phi(\underline{x}_i^T \underline{\beta})$

example: $Y_i =$ subject i has high blood pressure = $\begin{cases} 1 & \text{hypertension} \\ 0 & \text{no hypertension} \end{cases}$
we can only observe this binary version, but actually there is an underlying continuous r.v. (that we do not have) $Y_i^* =$ blood pressure
Indeed we can assume that Y_i is obtained starting from Y_i^* as

$$Y_i = \begin{cases} 1 & \text{if } Y_i^* > k \\ 0 & \text{if } Y_i^* \leq k \end{cases}$$

"subject i is considered to have high blood pressure if their pressure is above a threshold k "

→ For simplicity, we consider $k=0$ (it is sufficient to consider $Y_i^* - k$ for $k \neq 0$)
We assume a GAUSSIAN LINEAR MODEL on the LATENT VARIABLE Y_i^*

$$\left. \begin{aligned} Y_i^* &= \underline{x}_i^T \underline{\beta} + \epsilon_i \quad i=1, \dots, n \\ \epsilon_i &\text{ iid with distribution } \epsilon_i \sim N(0,1) \end{aligned} \right\} \Rightarrow Y_i^* \sim N(\underline{x}_i^T \underline{\beta}, 1) \text{ independent}$$

known variance = 1

However, we do not have Y_i^* , but only its dichotomized version Y_i :

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

what is $\mathbb{P}(Y_i=1) = \pi_i$?

$$\begin{aligned} \mathbb{P}(Y_i=1) &= \mathbb{P}(Y_i^* > 0) = 1 - \mathbb{P}(Y_i^* \leq 0) = 1 - \mathbb{P}(\underline{x}_i^T \underline{\beta} + \epsilon_i \leq 0) = \\ &= 1 - \mathbb{P}(\epsilon_i \leq -\underline{x}_i^T \underline{\beta}) \quad \epsilon_i \sim N(0,1) \\ &= 1 - \Phi(-\underline{x}_i^T \underline{\beta}) \\ &= \Phi(\underline{x}_i^T \underline{\beta}) \end{aligned}$$

⇒ $\pi_i = \Phi(\underline{x}_i^T \underline{\beta})$

which is exactly the model we assumed for Y_i (GLM).

Probit regression can be interpreted as a "simplification" of a Gaussian linear model, where we do not have all information on Y_i^* but only a dichotomized version.

