

Tutorial on Multinomial Logistic Regression

Javier R. Movellan

June 19, 2013

1 General Model

The inputs are n -dimensional vectors the outputs are c -dimensional vectors. The training sample consist of m input output pairs. We organize the example inputs as an $m \times n$ matrix x . The corresponding example outputs are organized as a $m \times c$ matrix y . The models under consideration make predictions

$$\hat{y} = h(u) \tag{1}$$

$$u = x\theta \tag{2}$$

where θ is a $n \times c$ weight matrix. Note θ_{ij} can be seen as the connection strength from input variable X_i to output variable U_j . We evaluate the

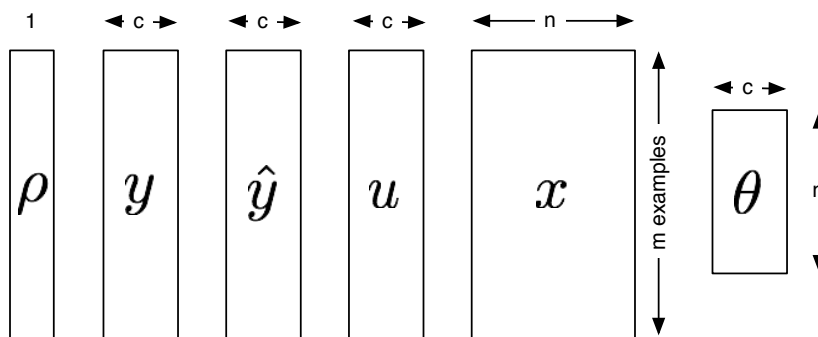


Figure 1: *There are m examples, n input variables and c output variables.*

optimality of \hat{y} , and thus of θ , using the following criterion

$$L(\theta) = \Phi(\theta) + \Gamma(\theta) \tag{3}$$

$$\Phi(\theta) = \sum_{j=1}^m w_j \rho_j \tag{4}$$

where $\Gamma(\theta)$ is a prior penalty function of the parameters,

$$\rho_j = f(y_{j\cdot}, \hat{y}_{j\cdot}) \tag{5}$$

measures the mismatch between the j^{th} rows of y and \hat{y} . The terms w_1, \dots, w_m are positive weights that capture the relative importance of each of the m input-output pairs.

Our goal is to find a matrix θ that minimizes L . A standard approach for minimizing L is the Newton-Raphson algorithm which calls for computation of the gradient and the Hessian matrix.

1.1 Gradient Vector

Let $\theta_{.i} \in \mathbb{R}^n$ be the i^{th} column of θ , i.e., the set of connection strengths from the n input units to the i^{th} output unit. The overall gradient of Φ with respect to θ is as follows

$$\nabla_{\text{vec}[\theta]}\Phi = \begin{pmatrix} \nabla_{\theta_{.1}}\Phi \\ \nabla_{\theta_{.2}}\Phi \\ \vdots \\ \nabla_{\theta_{.c}}\Phi \end{pmatrix} \quad (6)$$

where $\text{vec}[\theta]$ is the vectorized version of θ and the gradient of a vector q with respect to a vector p is defined as follows

$$\nabla_p q = \frac{\partial q}{\partial p} \quad (7)$$

Using the chain rule for gradients we get

$$\frac{\partial \Phi'}{\partial \theta_{.i}} = \frac{\partial u'_{.i}}{\partial \theta_{.i}} \frac{\partial \rho'}{\partial u_{.i}} \frac{\partial \Phi'}{\partial \rho} \quad (8)$$

Note

$$u_{.i} = (x\theta)_{.i} = x\theta_{.i} \quad (9)$$

Thus

$$\frac{\partial u'_{.i}}{\partial \theta_{.i}} = \frac{\partial \theta'_{.i} x'}{\partial \theta_{.i}} = x' \quad (10)$$

Moreover

$$\frac{\partial \rho'}{\partial u_{.i}} = \begin{pmatrix} \frac{\partial \rho_1}{\partial u_{1i}} & \cdots & \frac{\partial \rho_m}{\partial u_{1i}} \\ \vdots & \vdots & \vdots \\ \frac{\partial \rho_1}{\partial u_{mi}} & \cdots & \frac{\partial \rho_m}{\partial u_{mi}} \end{pmatrix} = \begin{pmatrix} \frac{\partial \rho_1}{\partial u_{1i}} & 0 & \cdots & 0 \\ 0 & \frac{\partial \rho_2}{\partial u_{2i}} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \frac{\partial \rho_m}{\partial u_{mi}} \end{pmatrix} \quad (11)$$

and

$$\frac{\partial \Phi'}{\partial \rho} = \begin{pmatrix} w_1 \\ \vdots \\ w_m \end{pmatrix} \quad (12)$$

Thus

$$\frac{\partial \Phi'}{\partial \theta_{.i}} = x' \begin{pmatrix} \frac{\partial \rho_1}{\partial u_{1i}} & 0 & \cdots & 0 \\ 0 & \frac{\partial \rho_2}{\partial u_{2i}} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \frac{\partial \rho_m}{\partial u_{mi}} \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{pmatrix} \quad (13)$$

Equivalently

$$\nabla_{\theta_{.i}} \Phi = \frac{\partial \Phi'}{\partial \theta_{.i}} = x' w \Psi_i \quad (14)$$

where

$$w = \begin{pmatrix} w_1 & 0 & \cdots & 0 \\ 0 & w_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & w_m \end{pmatrix} \quad (15)$$

and

$$\Psi_i = \begin{pmatrix} \frac{\partial \rho_1}{\partial u_{1i}} \\ \frac{\partial \rho_2}{\partial u_{2i}} \\ \vdots \\ \frac{\partial \rho_m}{\partial u_{mi}} \end{pmatrix} \quad (16)$$

Thus

$$\nabla_{\text{vec}[\theta]} \Phi = \begin{pmatrix} x' w \Psi_1 \\ \vdots \\ x' w \Psi_c \end{pmatrix} \quad (17)$$

1.2 Hessian Matrix

The Hessian of a scalar v with respect to a vector u is defined as follows

$$\nabla_u^2 v = \nabla_u (\nabla_u v) = \frac{\partial}{\partial u} (\nabla_u v)' \quad (18)$$

Thus

$$\nabla_{\text{vec}[\theta]}^2 \Phi = \nabla_{\text{vec}[\theta]} \begin{pmatrix} x'w\Psi_1 \\ \vdots \\ x'w\Psi_c \end{pmatrix} \quad (19)$$

$$= (\nabla_{\text{vec}[\theta]} x'w\Psi_1, \dots, x'w\Psi_c) \quad (20)$$

where

$$\nabla_{\text{vec}[\theta]} x'w\Psi_i = \nabla_{\text{vec}[\theta]} \Psi_i \nabla_{\Psi_i} x'w\Psi_i (\nabla_{\text{vec}[\theta]} \Psi_i) wx \quad (21)$$

Thus

$$\nabla_{\text{vec}[\theta]}^2 \Phi = \left((\nabla_{\text{vec}[\theta]} \Psi_1) wx, \dots, (\nabla_{\text{vec}[\theta]} \Psi_c) wx \right) \quad (22)$$

$$\nabla_{\text{vec}[\theta]}^2 \Phi = \left(\begin{pmatrix} \nabla_{\theta_1} \Psi_1 \\ \vdots \\ \nabla_{\theta_c} \Psi_1 \end{pmatrix} wx, \dots, \begin{pmatrix} \nabla_{\theta_1} \Psi_c \\ \vdots \\ \nabla_{\theta_c} \Psi_c \end{pmatrix} wx \right) \quad (23)$$

$$= \begin{pmatrix} (\nabla_{\theta_1} \Psi_1) wx & & (\nabla_{\theta_1} \Psi_c) wx \\ \vdots & \dots & \vdots \\ (\nabla_{\theta_c} \Psi_1) wx & & (\nabla_{\theta_c} \Psi_c) wx \end{pmatrix} \quad (24)$$

Note

$$\nabla_{\theta_i} \Psi_j = \nabla_{\theta_i} u_j \nabla_{u_j} \Psi_j = x' \Lambda_{ij} \quad (25)$$

where

$$\Lambda_{ij} = \nabla_{u_j} \Psi_j = \begin{pmatrix} \frac{\partial \Psi_{1j}}{\partial u_{1i}} & 0 & \dots & 0 \\ 0 & \frac{\partial \Psi_{2j}}{\partial u_{2i}} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \frac{\partial \Psi_{mj}}{\partial u_{mi}} \end{pmatrix} \quad (26)$$

Thus

$$\nabla_{\theta \cdot i} \Psi_j = \frac{\partial \Psi'_j}{\partial \theta \cdot i} = \begin{pmatrix} x_{11} \frac{\partial \Psi_{1j}}{\partial u_{1i}} & \cdots & x_{m1} \frac{\partial \Psi_{mj}}{\partial u_{mi}} \\ \vdots & \vdots & \vdots \\ x_{1n} \frac{\partial \Psi_{1j}}{\partial u_{ni}} & \cdots & x_{mn} \frac{\partial \Psi_{mj}}{\partial u_{mi}} \end{pmatrix} = x' \Lambda_{ij} \quad (27)$$

where

$$\Lambda_{ij} = \begin{pmatrix} \frac{\partial \Psi_{1j}}{\partial u_{1i}} & 0 & \cdots & 0 \\ 0 & \frac{\partial \Psi_{2j}}{\partial u_{2i}} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & \frac{\partial \Psi_{mj}}{\partial u_{mi}} \end{pmatrix} \quad (28)$$

Thus

$$\nabla_{\text{vec}[\theta]}^2 \Phi = \begin{pmatrix} x' \Lambda_{11} w x & \cdots & x' \Lambda_{1c} w x \\ \vdots & \cdots & \vdots \\ x' \Lambda_{c1} w x & \cdots & x' \Lambda_{cc} w x \end{pmatrix} \quad (29)$$

2 Quadratic Priors

Let

$$\Gamma(\theta) = -\frac{1}{2} (\text{vec}[\theta] - \text{vec}[\mu])' \sigma^{-1} (\text{vec}[\theta] - \text{vec}[\mu]) \quad (30)$$

then

$$\nabla_{\text{vec}[\theta]} \Gamma(\theta) = -\sigma^{-1} (\text{vec}[\theta] - \text{vec}[\mu]) \quad (31)$$

$$\nabla_{\text{vec}[\theta]}^2 \Gamma = \nabla_{\text{vec}[\theta]} \nabla_{\text{vec}[\theta]} \Gamma = -\sigma^{-1} \quad (32)$$

3 Newton-Raphson Optimization

$$\text{vec}[\theta^{(k+1)}] = \text{vec}[\theta^{(k)}] - (\nabla_{\text{vec}[\theta]}^2 L(\theta^{(k)}))^{-1} \nabla_{\text{vec}[\theta]} L(\theta^{(k)}) \quad (33)$$

where $\theta^{(k)}$ is the value of θ after k iterations of the optimization algorithm.

4 Multivariate Linear Regression

In this case

$$\hat{y}_{i\cdot} = u_i. \quad (34)$$

$$\rho_i = \sum_k (\hat{y}_{ik} - y_{ik})^2 \quad (35)$$

Thus

$$\Psi_i = \hat{y}_{i\cdot} - y_{i\cdot} \quad (36)$$

$$\Lambda_{ij} = I_m \quad (37)$$

where I_m is the $m \times m$ identity matrix. Thus

$$\frac{\partial}{\partial \theta_{\cdot i}} \frac{\partial \Phi}{\partial \theta'_{\cdot j}} = x' w x \quad (38)$$

5 Multinomial Logistic Regression

Let

$$\hat{y}_{ij} = f_j(u_i) = \frac{e^{u_{ij}}}{\sum_{k=1}^c e^{u_{ik}}} \quad (39)$$

and ρ_i the negative log-likelihood of the output vector y_i given the input vector x_i .

$$\rho_i = - \sum_{k=1}^c y_{ik} \log \hat{y}_{ik} \quad (40)$$

Note

$$\frac{\partial \hat{y}_{ij}}{\partial u_{ik}} = \hat{y}_{ij} \frac{\partial \log \hat{y}_{ij}}{\partial u_{ik}} = \hat{y}_{ij} (\delta_{jk} - \hat{y}_{ik}) \quad (41)$$

Thus

$$\frac{\partial \rho_i}{\partial u_{ik}} = - \sum_{j=1}^c \frac{y_{ij}}{\hat{y}_{ij}} \frac{\partial \hat{y}_{ij}}{\partial u_{ik}} = \sum_{j=1}^c y_{ij} (\delta_{jk} - \hat{y}_{ik}) = \hat{y}_{ik} - y_{ik} \quad (42)$$

where we used the fact that $\sum_j y_{ij} = 1$. Thus

$$\Psi_i = \begin{pmatrix} \frac{\partial \rho_1}{\partial u_{1i}} \\ \frac{\partial \rho_2}{\partial u_{2i}} \\ \vdots \\ \frac{\partial \rho_m}{\partial u_{mi}} \end{pmatrix} = \begin{pmatrix} y_{1i} - \hat{y}_{1i} \\ y_{2i} - \hat{y}_{2i} \\ \vdots \\ y_{mi} - \hat{y}_{mi} \end{pmatrix} = \hat{y}_{\cdot i} - y_{\cdot i} \quad (43)$$

Moreover

$$\Lambda_{ij} = \begin{pmatrix} \frac{\partial \Psi_{1j}}{\partial u_{1i}} & 0 & \dots & 0 \\ 0 & \frac{\partial \Psi_{2j}}{\partial u_{2i}} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \frac{\partial \Psi_{mj}}{\partial u_{mi}} \end{pmatrix} \quad (44)$$

where

$$\frac{\partial \Psi_{kj}}{\partial u_{ki}} = \frac{\partial \hat{y}_{kj}}{\partial u_{ki}} = \hat{y}_{kj}(\delta_{ij} - \hat{y}_{kj}) \quad (45)$$

5.1 Relationship to Linear Regression

Note that the gradient in multinomial logistic regression is identical to the gradient in multivariate linear regression.

$$\nabla_{\theta_{\cdot i}} \Phi = \hat{y}_{\cdot i} - y_{\cdot i} \quad (46)$$

The Hessians would also be very similar. In linear regression

$$\frac{\partial}{\partial \theta_{\cdot i}} \frac{\partial \Phi}{\partial \theta'_{\cdot j}} = x' w x \quad (47)$$

and in logistic regression

$$\frac{\partial}{\partial \theta_{\cdot i}} \frac{\partial \Phi}{\partial \theta'_{\cdot j}} = x' w \Lambda_{ij} x \quad (48)$$

which can be seen as special case of linear regression where the weight matrix w is substituted by the $w \Lambda_{ij}$ matrix.

5.2 Summary

- **Training Data**

- **Inputs** $x \in \mathbb{R}^m \times \mathbb{R}^n$. Rows are examples, columns are input variables.
- **Outputs** $y \in \mathbb{R}^m \times \mathbb{R}^c$. Rows are examples, columns are labels or label probabilities. All entries are non-negative and each row add up to 1.
- **Example Weights** $w \in \mathbb{R}^m \times \mathbb{R}^r$. Diagonal matrix. Each term is positive and represents the relative importance of each example.

- **Gradients**

$$\nabla_{\text{vec}[\theta]} L = \begin{pmatrix} \nabla_{\theta_1} \Phi \\ \nabla_{\theta_2} \Phi \\ \vdots \\ \nabla_{\theta_c} \Phi \end{pmatrix} + \nabla_{\text{vec}[\theta]} \Gamma \quad (49)$$

where

$$\nabla_{\theta_i} \Phi = \frac{\partial \Phi'}{\partial \theta_i} = x' w \Psi_i \quad (50)$$

and

$$\Psi_i = \begin{pmatrix} \frac{\partial \rho_1}{\partial u_{1i}} \\ \frac{\partial \rho_2}{\partial u_{2i}} \\ \vdots \\ \frac{\partial \rho_m}{\partial u_{mi}} \end{pmatrix} = \hat{y}_{\cdot i} - y_{\cdot i} \quad (51)$$

and

$$\nabla_{\text{vec}[\theta]} \Gamma = -\sigma^{-1}(\text{vec}[\theta] - \text{vec}[\mu]) \quad (52)$$

- **Hessian Matrices**

$$\nabla_{\text{vec}[\theta]}^2 L = \nabla_{\text{vec}[\theta]}^2 \Phi + \nabla_{\text{vec}[\theta]}^2 \Gamma \quad (53)$$

$$\nabla_{\text{vec}[\theta]}^2 \Phi = \begin{pmatrix} \frac{\partial}{\partial \theta_{.1}} \frac{\partial \Phi}{\partial \theta'_{.1}} & \cdots & \frac{\partial}{\partial \theta_{.1}} \frac{\partial \Phi}{\partial \theta'_{.c}} \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial \theta_{.c}} \frac{\partial \Phi}{\partial \theta'_{.1}} & \cdots & \frac{\partial}{\partial \theta_{.c}} \frac{\partial \Phi}{\partial \theta'_{.c}} \end{pmatrix} \quad (54)$$

where

$$\frac{\partial}{\partial \theta_{.i}} \frac{\partial \Phi}{\partial \theta'_{.j}} = x' \Lambda_{ij} w x \quad (55)$$

where

$$\Lambda_{ij} = \begin{pmatrix} \frac{\partial \Psi_{1j}}{\partial u_{1i}} & 0 & \cdots & 0 \\ 0 & \frac{\partial \Psi_{2j}}{\partial u_{2i}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\partial \Psi_{mj}}{\partial u_{mi}} \end{pmatrix} \quad (56)$$

and

$$\frac{\partial \Psi_{kj}}{\partial u_{ki}} = \frac{\partial \hat{y}_{kj}}{\partial u_{ki}} = \hat{y}_{kj} (\delta_{ij} - \hat{y}_{kj}) \quad (57)$$

and

$$\nabla_{\text{vec}[\theta]}^2 \Gamma = -\sigma^{-1} \quad (58)$$

- **Learning Rule (Newton-Raphson)**

$$\text{vec}[\theta^{(k+1)}] = \text{vec}[\theta^{(k)}] - (\nabla_{\text{vec}[\theta]}^2 L(\theta^{(k)})^{-1} \nabla_{\text{vec}[\theta]} L(\theta^{(k)})) \quad (59)$$

where $\theta^{(k)}$ is the value of θ after k iterations of the learning rule.

6 Appendix: L-p priors

From a Bayesian point of view the Φ function can be seen as a log-likelihood term. At times it may be useful to add a log prior term over θ , also known as a regularization term. Most prior terms of interest have the following form

$$\Gamma = \alpha \sum_{i=1}^n \sum_{j=1}^c \frac{1}{p} \left(h(\theta_{i,j}) \right)^p \quad (60)$$

where α is a non-negative constant. If h is the absolute value function then Γ is the L - p norm of θ . Two popular options are $p = 2$ and $p = 1$. Note the absolute value function is not differentiable, and thus when p is odd, it useful to approximate it with a differentiable function. Here we will use the log hyperbolic cosine function.

$$h(x) = \frac{1}{\beta} \log(2\cosh(\beta x)) \quad (61)$$

where $\beta \geq 0$ is a gain parameter and

$$\cosh(x) = \frac{e^x + e^{-x}}{2} \quad (62)$$

Note

$$\lim_{\beta \rightarrow \infty} h(x) = \text{abs}(x) \quad (63)$$

$$h'(x) = \frac{dh(x)}{dx} = \frac{e^{\beta x} - e^{-\beta x}}{e^{\beta x} + e^{-\beta x}} = \tanh(x) \quad (64)$$

$$\lim_{\beta \rightarrow \infty} \frac{dh(x)}{dx} = \text{sign}(x) \quad (65)$$

and

$$h''(x) = \frac{d^2h(x)}{dx^2} = \frac{d \tanh(\beta x)}{dx} = \beta (1 - \tanh^2(\beta x)) \quad (66)$$

6.1 Gradient and Hessian

$$\frac{\partial \Gamma}{\partial \theta_{ij}} = \alpha \left(h^{p-1}(\theta_{ij}) h'(\theta_{ij}) \right) \quad (67)$$

$$= \alpha \left(\left(\frac{1}{\beta} \log(2\cosh(\beta x)) \right)^{p-1} \tanh(x) \right) \quad (68)$$

The $nc \times nc$ Hessian matrix is diagonal. Each term in the diagonal has the following form

$$\frac{\partial^2 \Gamma}{\partial \theta_{ij}^2} = \alpha \left((p-1) h(\theta_{ij})^{p-2} (h'(\theta_{ij}))^2 + h(\theta_{ij})^{p-1} h''(\theta_{ij}) \right) \quad (69)$$

Thus

$$\frac{\partial^2 \Gamma}{\partial \theta_{ij}^2} = \alpha \left((p-1) \left(\frac{1}{\beta} \log(2 \cosh(\beta x)) \right)^{p-2} (\tanh(\beta x))^2 \right) \quad (70)$$

$$+ \left(\frac{1}{\beta} \log(2 \cosh(\beta x)) \right)^{p-1} \beta (1 - \tanh^2(\beta x)) \quad (71)$$

Note for $p = 1$, we get

$$\frac{\partial \Gamma}{\partial \theta_{ij}} = \alpha \frac{1}{\beta} \tanh(\beta \theta_{ij}) \quad (72)$$

$$\frac{\partial^2 \Gamma}{\partial \theta_{ij}^2} = \alpha \beta (1 - \tanh^2(\beta \theta_{ij})) \quad (73)$$

$$(74)$$

and for $p = 1$, and $\beta \rightarrow \infty$ we get

$$\frac{\partial \Gamma}{\partial \theta_{ij}} = \alpha \theta_{ij} \quad (75)$$

$$\frac{\partial^2 \Gamma}{\partial \theta_{ij}^2} = \alpha \quad (76)$$

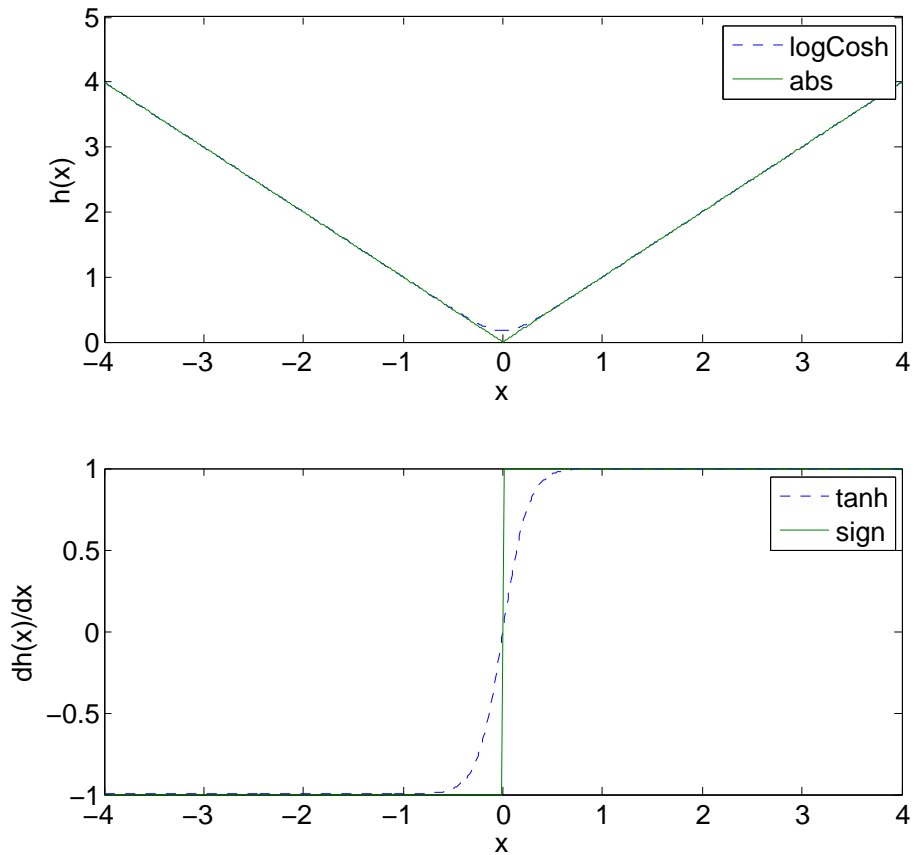


Figure 2: *The absolute value function can be approximated by the log of a hyperbolic cosine. The derivative of the absolute value is the sign function. It can be approximated by the derivative of the log cosh function, which is the hyperbolic tangent. In this figure we used gain $\beta = 4$.*