

Lecture 05 – Artificial neural networks deep learning

Prof. João Fernando Mari

joaofmari.github.io

joaof.mari@ufv.br

- Multilayer perceptron
- Network architecture
- Backpropagation – Algorithm and equations
- Backpropagation – Functions and derivatives
- Backpropagation – Examples

MULTILAYER PERCEPTRON

Multilayer perceptron

NATURE VOL. 323 9 OCTOBER 1986

LETTERS TO NATURE

523

Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

There have been many attempts to design self-organizing neural networks. The aim is to find a powerful synaptic modification rule that will allow an arbitrarily connected neural network to develop an internal structure that is appropriate for a particular task domain. The task is specified by giving the desired state vector of the output units for each state vector of the input units. If the input units are directly connected to the output units it is relatively easy to find learning rules that iteratively adjust the relative strengths of the connections so as to progressively reduce the difference between the actual and desired output vectors². Learning becomes more interesting but

more difficult when we introduce hidden units whose actual or desired states are not specified by the task. (In perceptrons, there are 'feature analysers' between the input and output that are not true hidden units because their input connections are fixed by hand, so their states are completely determined by the input vector: they do not learn representations.) The learning procedure must decide under what circumstances the hidden units should be active in order to help achieve the desired input-output behaviour. This amounts to deciding what these units should represent. We demonstrate that a general purpose and relatively simple procedure is powerful enough to construct appropriate internal representations.

The simplest form of the learning procedure is for layered networks which have a layer of input units at the bottom, any number of intermediate layers, and a layer of output units at the top. Connections within a layer or from higher to lower layers are forbidden, but connections can skip intermediate layers. An input vector is presented to the network by setting the states of the input units. Then the states of the units in each layer are determined by applying equations (1) and (2) to the connections coming from lower layers. All units within a layer have their states set in parallel, but different layers have their states set sequentially, starting at the bottom and working upwards until the states of the output units are determined.

The total input, x_j , to unit j is a linear function of the outputs, y_i , of the units that are connected to j and of the weights, w_{ij} , on these connections

$$x_j = \sum_i y_i w_{ij} \quad (1)$$

Units can be given biases by introducing an extra input to each unit which always has a value of 1. The weight on this extra input is called the bias and is equivalent to a threshold of the opposite sign. It can be treated just like the other weights.

A unit has a real-valued output, y_j , which is a non-linear function of its total input

$$y_j = \frac{1}{1 + e^{-x_j}} \quad (2)$$

¹ To whom correspondence should be addressed

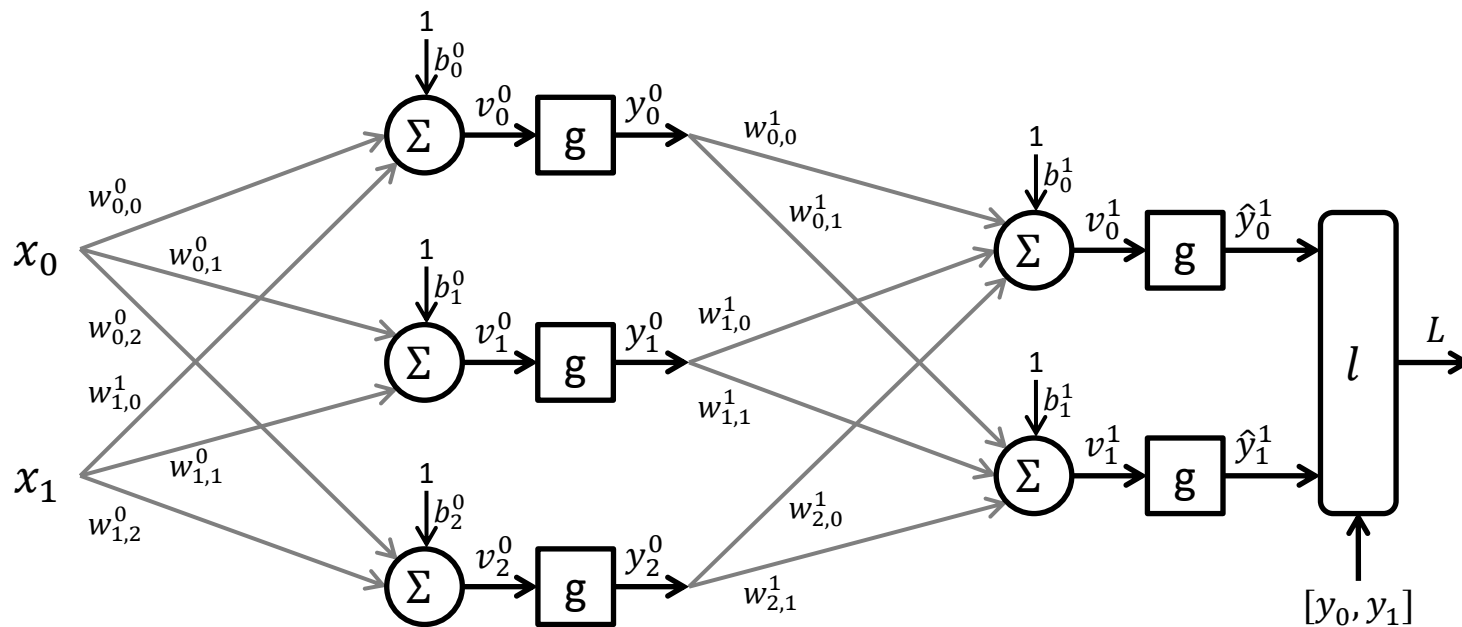


David E. Rumelhart
(1942 – 2011)
Cognitive psychology

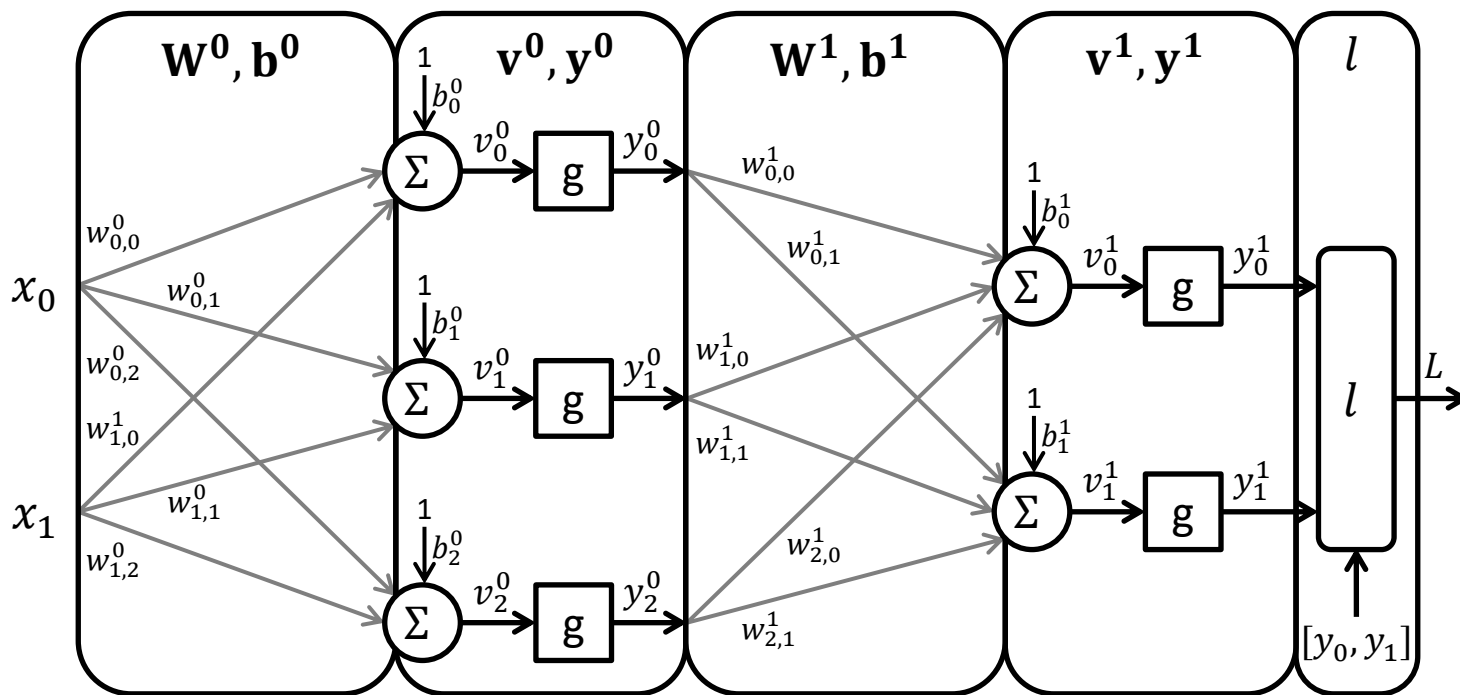


Geoffrey E. Hinton
(1947 -)
Cognitive psychology
Computer science

Multilayer perceptron

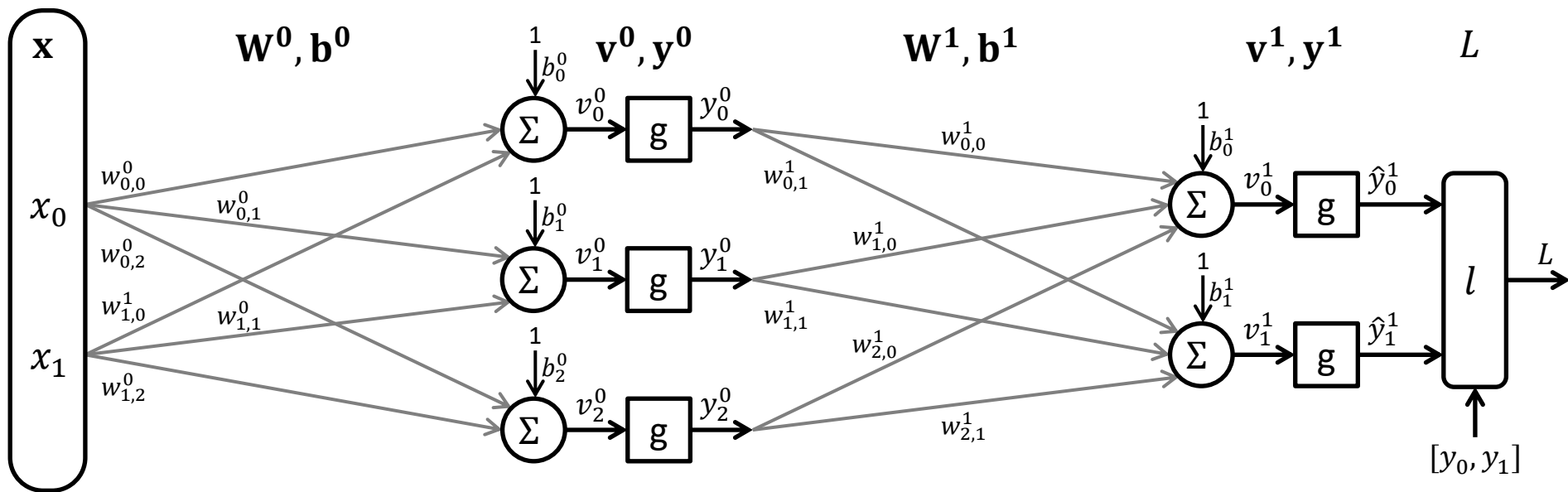


Multilayer perceptron



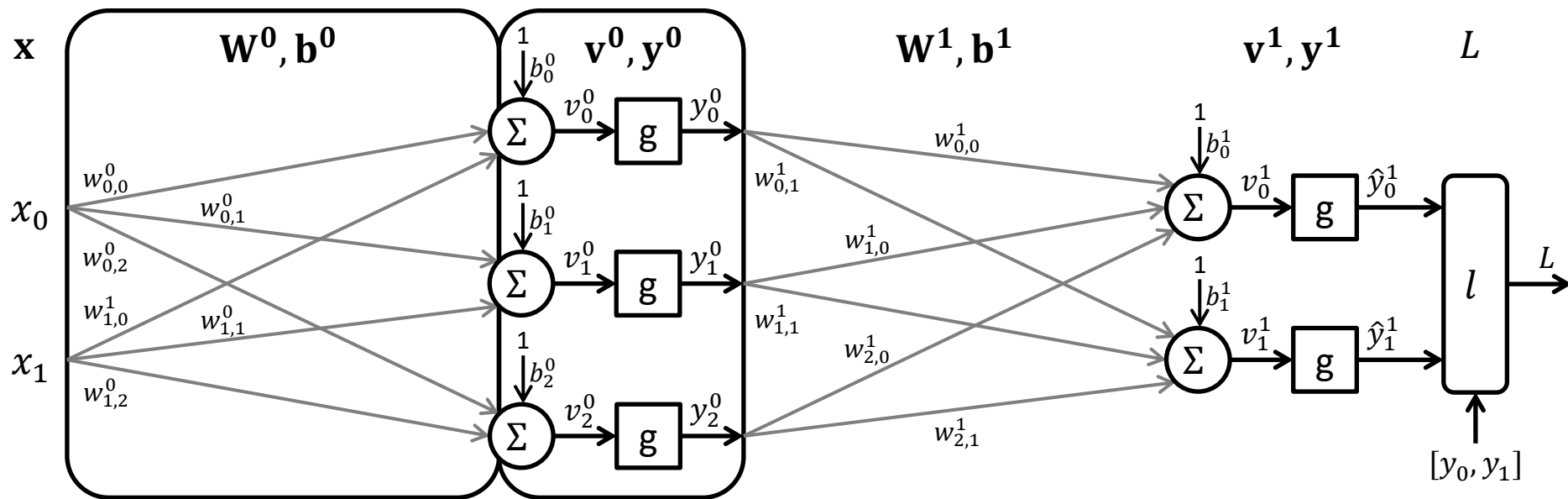
NETWORK ARCHITECTURE

Network architecture



$$\mathbf{x} = [x_0 \quad x_1]$$

Network architecture

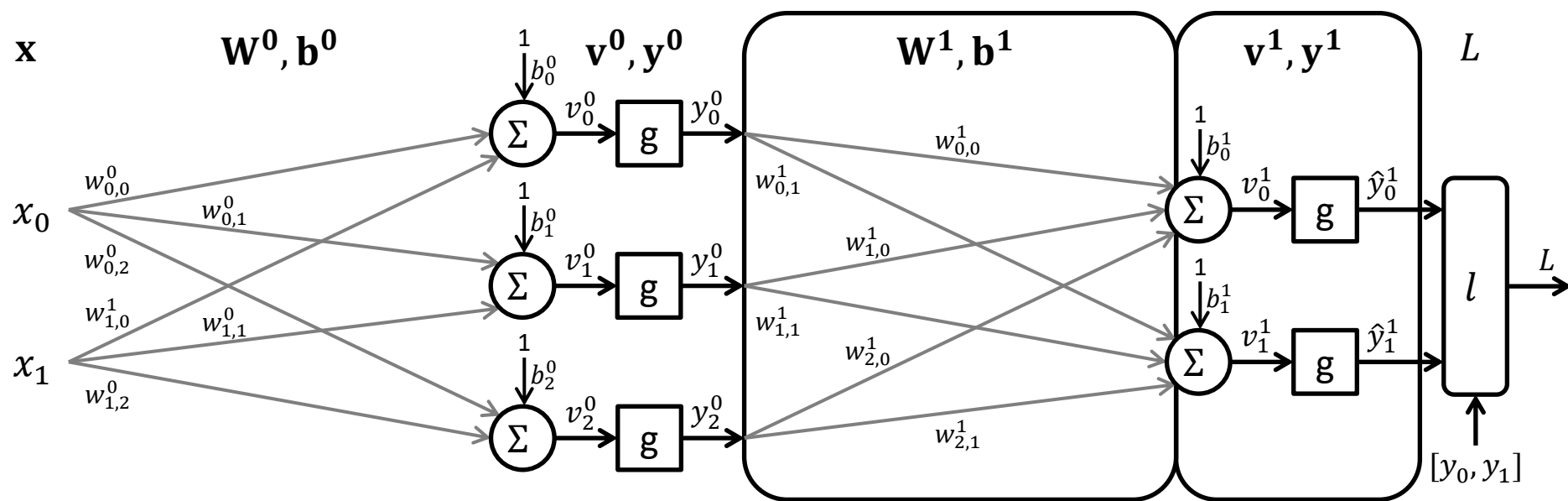


$$\mathbf{x} = [x_0 \quad x_1] \quad \mathbf{W}^0 = \begin{bmatrix} w_{0,0}^0 & w_{0,1}^0 & w_{0,2}^0 \\ w_{1,0}^0 & w_{1,1}^0 & w_{1,2}^0 \end{bmatrix} \quad \mathbf{v}^0 = [v_0^0 \quad v_1^0 \quad v_2^0]$$

$$\mathbf{y}^0 = [y_0^0 \quad y_1^0 \quad y_2^0]$$

$$\mathbf{b}^0 = [b_0^0 \quad b_1^0 \quad b_2^0]$$

Network architecture



$$\mathbf{x} = [x_0 \quad x_1]$$

$$\mathbf{W}^0 = \begin{bmatrix} w_{0,0}^0 & w_{0,1}^0 & w_{0,2}^0 \\ w_{1,0}^0 & w_{1,1}^0 & w_{1,2}^0 \end{bmatrix}$$

$$\mathbf{v}^0 = [v_0^0 \quad v_1^0 \quad v_2^0]$$

$$\mathbf{y}^0 = [y_0^0 \quad y_1^0 \quad y_2^0]$$

$$\mathbf{W}^1 = \begin{bmatrix} w_{0,0}^1 & w_{0,1}^1 \\ w_{1,0}^1 & w_{1,1}^1 \\ w_{2,0}^1 & w_{2,1}^1 \end{bmatrix}$$

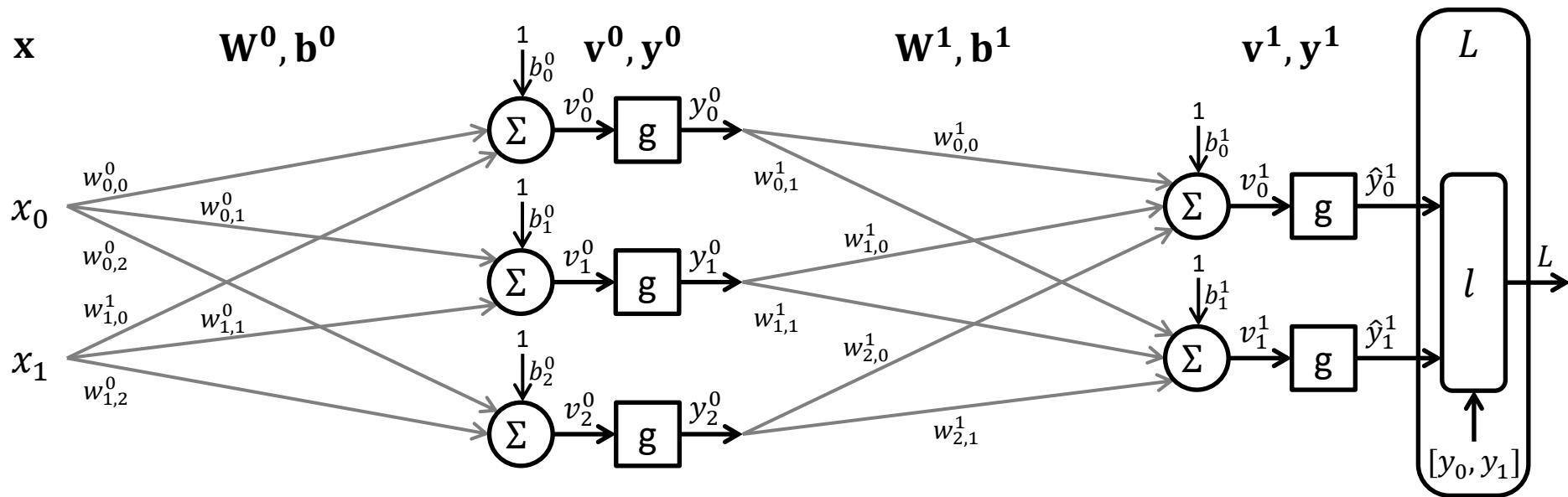
$$\mathbf{v}^1 = [v_0^1 \quad v_1^1]$$

$$\hat{\mathbf{y}}^1 = [\hat{y}_0^1 \quad \hat{y}_1^1]$$

$$\mathbf{b}^0 = [b_0^0 \quad b_1^0 \quad b_2^0]$$

$$\mathbf{b}^1 = [b_0^1 \quad b_1^1]$$

Network architecture



$$\mathbf{x} = [x_0 \quad x_1] \quad \mathbf{W}^0 = \begin{bmatrix} w_{0,0}^0 & w_{0,1}^0 & w_{0,2}^0 \\ w_{1,0}^0 & w_{1,1}^0 & w_{1,2}^0 \end{bmatrix}$$

$$\mathbf{b}^0 = [b_0^0 \quad b_1^0 \quad b_2^0]$$

$$\mathbf{v}^0 = [v_0^0 \quad v_1^0 \quad v_2^0] \\ \mathbf{y}^0 = [y_0^0 \quad y_1^0 \quad y_2^0]$$

$$\mathbf{W}^1 = \begin{bmatrix} w_{0,0}^1 & w_{0,1}^1 \\ w_{1,0}^1 & w_{1,1}^1 \\ w_{2,0}^1 & w_{2,1}^1 \end{bmatrix}$$

$$\mathbf{b}^1 = [b_0^1 \quad b_1^1]$$

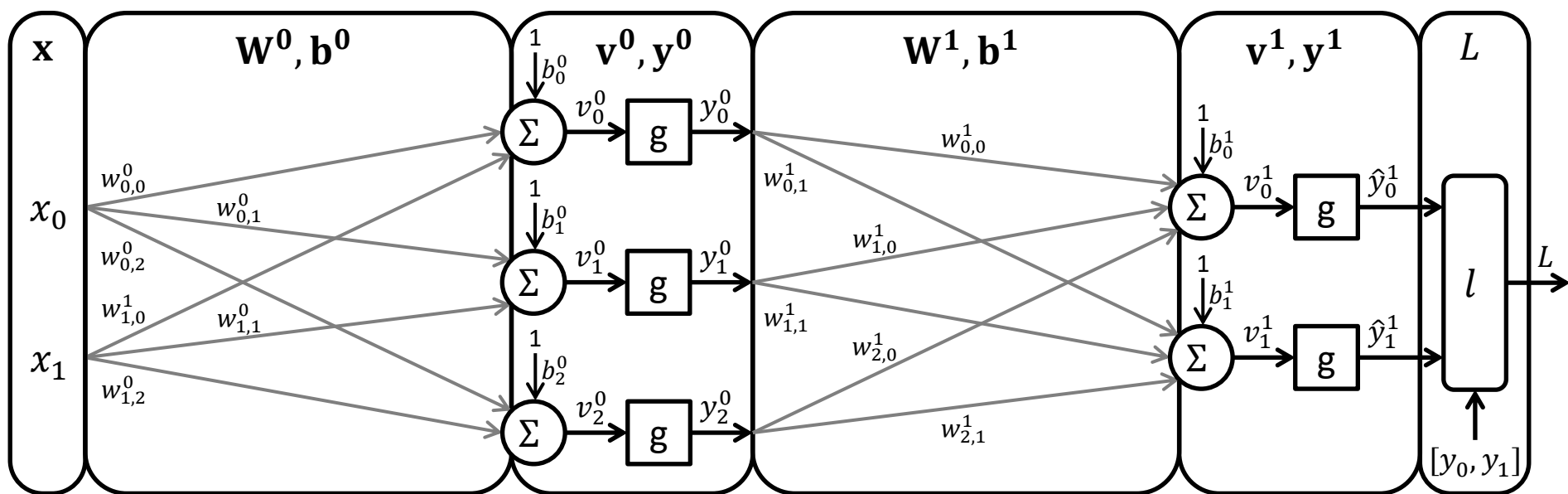
$$\mathbf{v}^1 = [v_0^1 \quad v_1^1] \\ \hat{\mathbf{y}}^1 = [\hat{y}_0^1 \quad \hat{y}_1^1]$$

$$\mathbf{y} = [y_0 \quad y_1]$$

$$\mathbf{l} = [l_0 \quad l_1]$$

$$L = l_0 + l_1$$

Network architecture



$$\mathbf{x} = [x_0 \quad x_1]$$

$$\mathbf{W}^0 = \begin{bmatrix} w_{0,0}^0 & w_{0,1}^0 & w_{0,2}^0 \\ w_{1,0}^0 & w_{1,1}^0 & w_{1,2}^0 \end{bmatrix}$$

$$\mathbf{b}^0 = [b_0^0 \quad b_1^0 \quad b_2^0]$$

$$\mathbf{v}^0 = [v_0^0 \quad v_1^0 \quad v_2^0]$$

$$\mathbf{y}^0 = [y_0^0 \quad y_1^0 \quad y_2^0]$$

$$\mathbf{W}^1 = \begin{bmatrix} w_{0,0}^1 & w_{0,1}^1 \\ w_{1,0}^1 & w_{1,1}^1 \\ w_{2,0}^1 & w_{2,1}^1 \end{bmatrix}$$

$$\mathbf{b}^1 = [b_0^1 \quad b_1^1]$$

$$\mathbf{v}^1 = [v_0^1 \quad v_1^1]$$

$$\hat{\mathbf{y}}^1 = [\hat{y}_0^1 \quad \hat{y}_1^1]$$

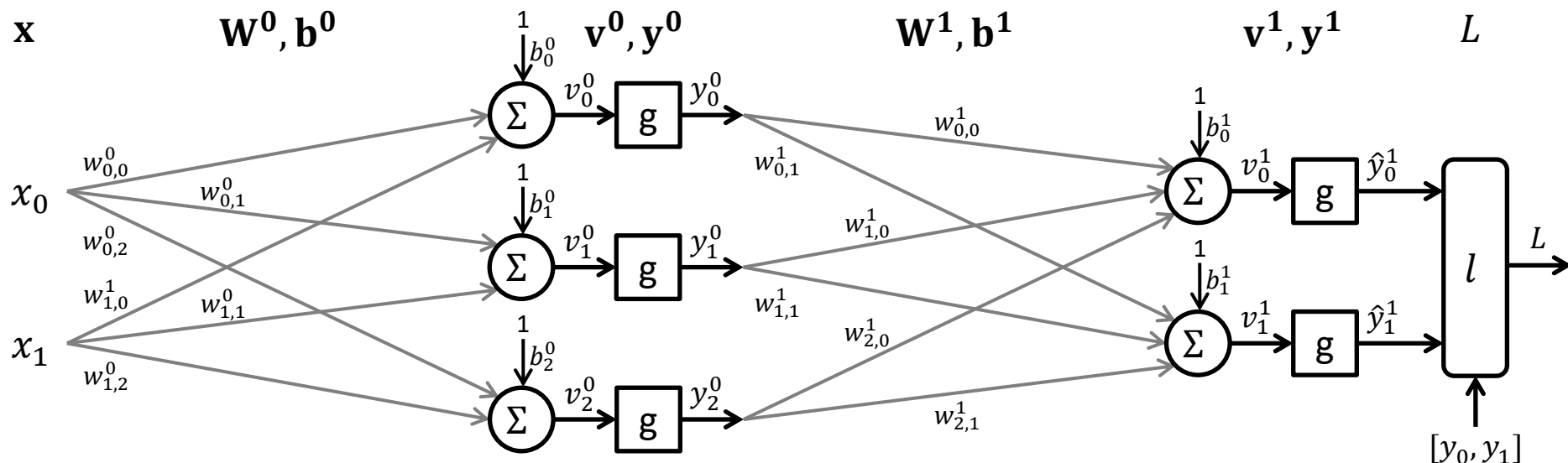
$$\mathbf{y} = [y_0 \quad y_1]$$

$$\mathbf{l} = [l_0 \quad l_1]$$

$$L = l_0 + l_1$$

BACKPROPAGATION – ALGORITHM AND EQUATIONS

Backpropagation – Algorithm and equations



$$\mathbf{v}^0 = \mathbf{x}\mathbf{W}^0 + \mathbf{b}^0$$

$$\mathbf{y}^0 = \mathbf{g}^0(\mathbf{v}^0)$$

$$\mathbf{v}^1 = \mathbf{y}^0\mathbf{W}^1 + \mathbf{b}^1$$

$$\mathbf{y}^1 = \mathbf{g}^1(\mathbf{v}^1)$$

$$l = l(\mathbf{y}, \hat{\mathbf{y}}^1)$$

$$L = l_0 + l_1$$

$$\mathbf{v}^0 = [x_0 \quad x_1] \begin{bmatrix} w_{0,0}^0 & w_{0,1}^0 & w_{0,2}^0 \\ w_{1,0}^0 & w_{1,1}^0 & w_{1,2}^0 \end{bmatrix} + [b_0^0 \quad b_1^0 \quad b_2^0]$$

$$\mathbf{v}^0 = [x_0 w_{0,0}^0 + x_1 w_{1,0}^0 + b_0^0 \quad x_0 w_{0,1}^0 + x_1 w_{1,1}^0 + b_1^0 \quad x_0 w_{0,2}^0 + x_1 w_{1,2}^0 + b_2^0]$$

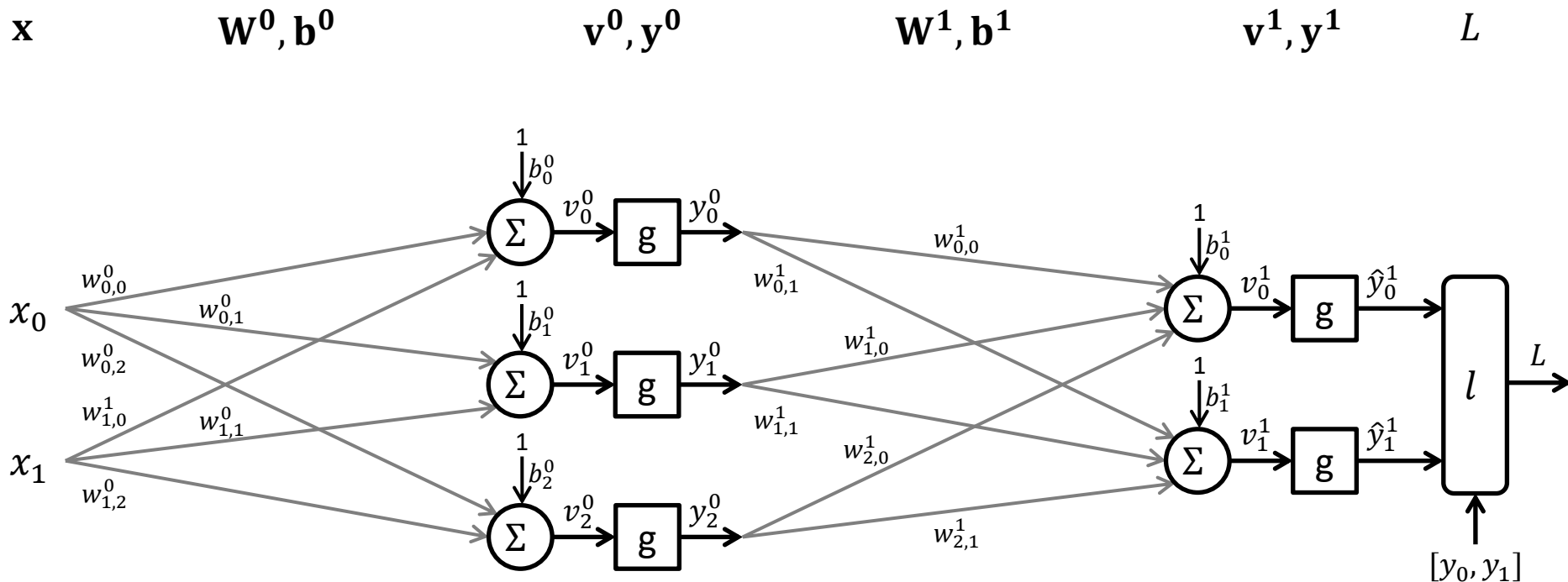
$$\mathbf{y}^0 = [g(v_0^0) \quad g(v_1^0) \quad g(v_2^0)]$$

$$\mathbf{v}^1 = [y_0^0 \quad y_1^0 \quad y_2^0] \begin{bmatrix} w_{0,0}^1 & w_{0,1}^1 \\ w_{1,0}^1 & w_{1,1}^1 \\ w_{2,0}^1 & w_{2,1}^1 \end{bmatrix} + [b_0^1 \quad b_1^1]$$

$$\mathbf{v}^1 = [y_0^0 w_{0,0}^1 + y_1^0 w_{1,0}^1 + y_2^0 w_{2,0}^1 + b_0^1 \quad y_0^0 w_{0,1}^1 + y_1^0 w_{1,1}^1 + y_2^0 w_{2,1}^1 + b_1^1]$$

$$\hat{\mathbf{y}}^1 = [g(v_0^1) \quad g(v_1^1)]$$

Backpropagation – Algorithm and equations



$$\frac{\partial L}{\partial \mathbf{W}^0} = \frac{\partial L}{\partial \mathbf{y}^0} \times \frac{\partial \mathbf{y}^0}{\partial \mathbf{v}^0} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{W}^0}$$

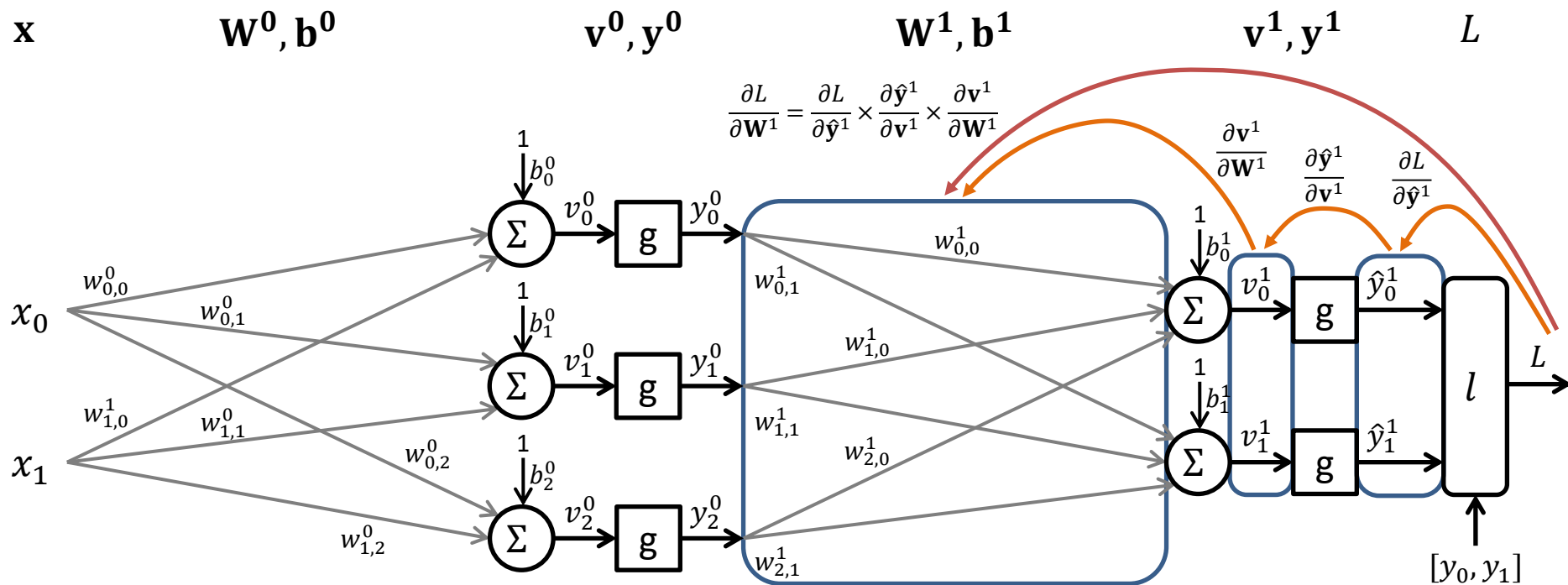
$$\frac{\partial L}{\partial \mathbf{y}^0} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{y}^0}$$

$$\frac{\partial L}{\partial \mathbf{b}^0} = \frac{\partial L}{\partial \mathbf{y}^0} \times \frac{\partial \mathbf{y}^0}{\partial \mathbf{v}^0} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{b}^0}$$

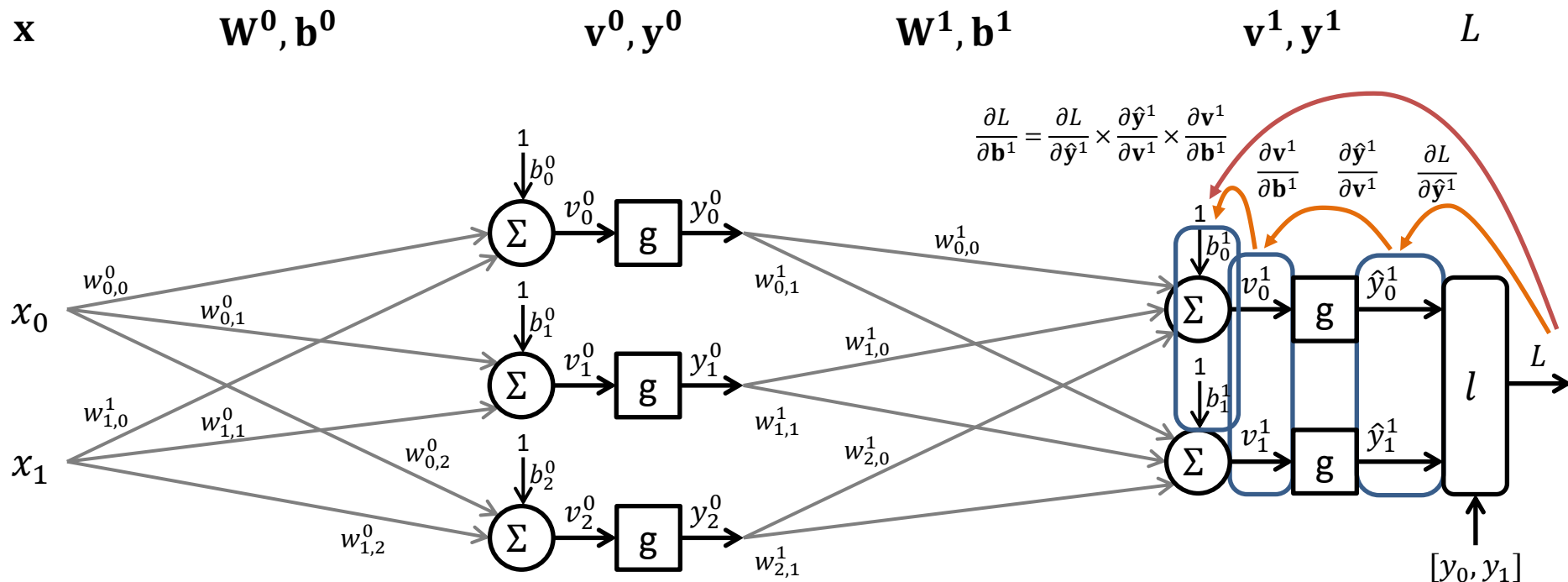
$$\frac{\partial L}{\partial \mathbf{W}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{W}^1}$$

$$\frac{\partial L}{\partial \mathbf{b}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{b}^1}$$

Backpropagation – Algorithm and equations

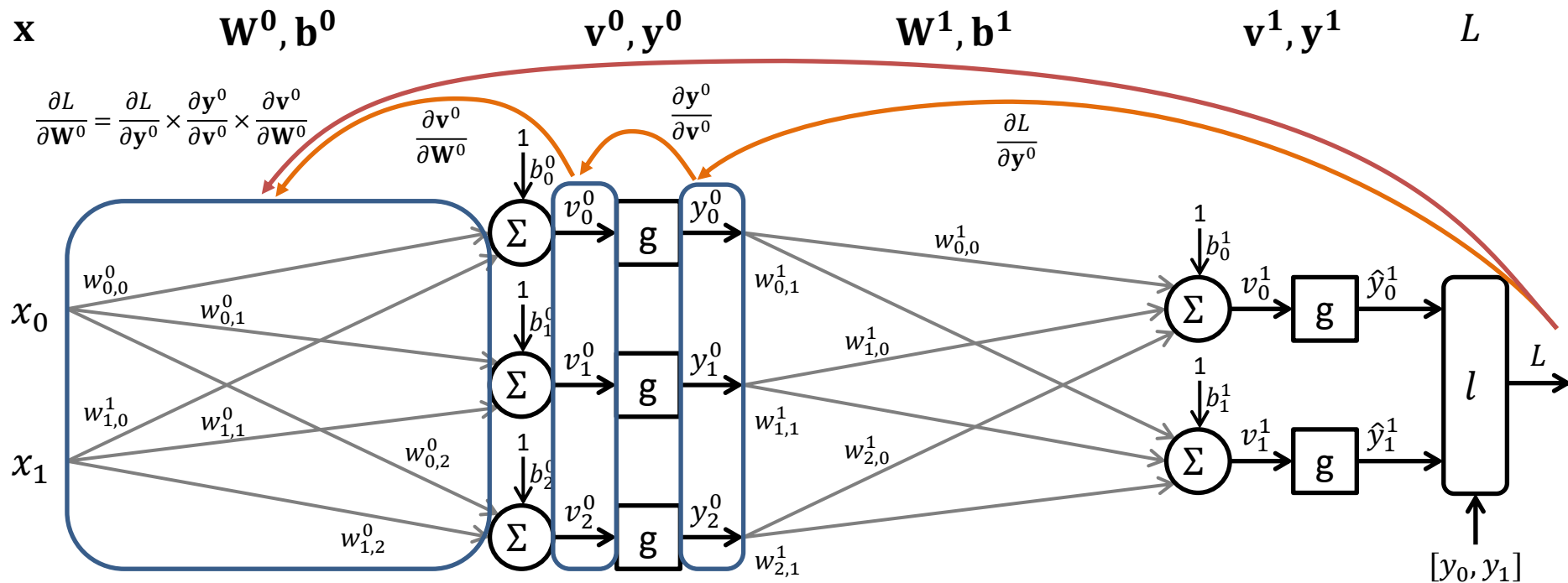


Backpropagation – Algorithm and equations



$$\frac{\partial L}{\partial \mathbf{b}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{b}^1}$$

Backpropagation – Algorithm and equations



$$\frac{\partial L}{\partial W^0} = \frac{\partial L}{\partial y^0} \times \frac{\partial y^0}{\partial v^0} \times \frac{\partial v^0}{\partial W^0}$$

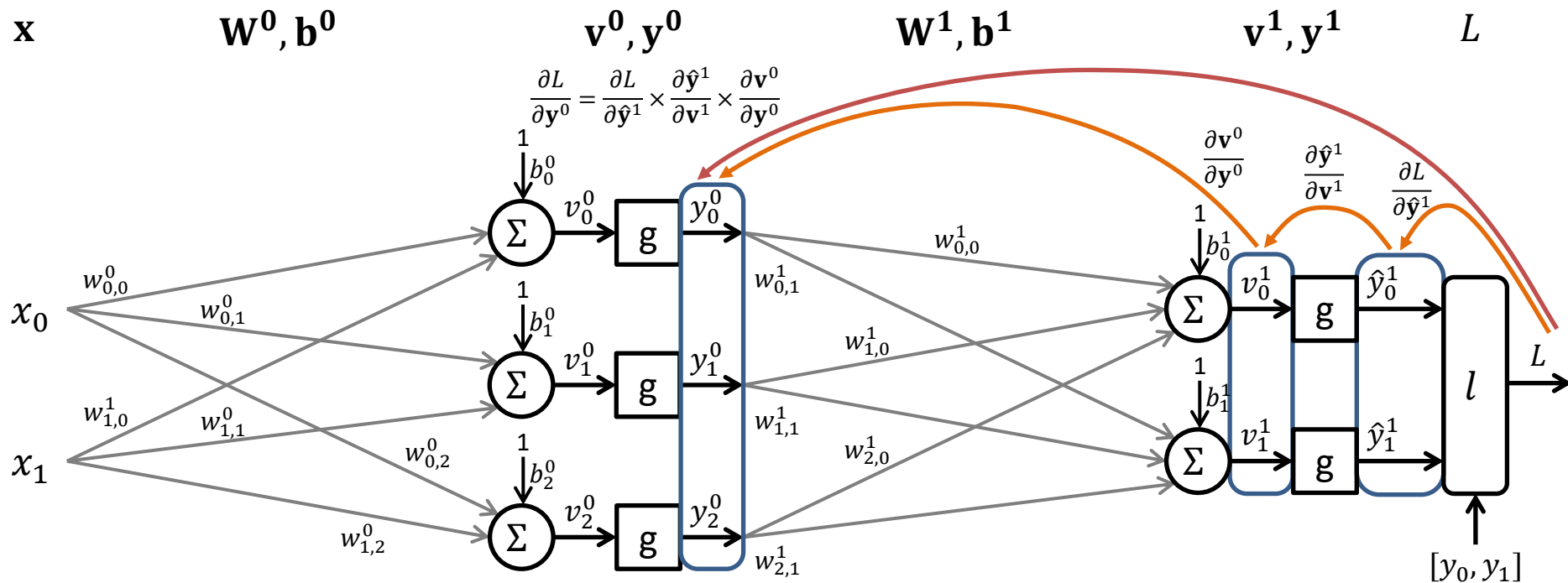
$$\frac{\partial L}{\partial y^0} = \frac{\partial L}{\partial \hat{y}^1} \times \frac{\partial \hat{y}^1}{\partial v^1} \times \frac{\partial v^1}{\partial y^0}$$

$$\frac{\partial L}{\partial b^0} = \frac{\partial L}{\partial y^0} \times \frac{\partial y^0}{\partial v^0} \times \frac{\partial v^0}{\partial b^0}$$

$$\frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial \hat{y}^1} \times \frac{\partial \hat{y}^1}{\partial v^1} \times \frac{\partial v^1}{\partial W^1}$$

$$\frac{\partial L}{\partial b^1} = \frac{\partial L}{\partial \hat{y}^1} \times \frac{\partial \hat{y}^1}{\partial v^1} \times \frac{\partial v^1}{\partial b^1}$$

Backpropagation – Algorithm and equations



$$\frac{\partial L}{\partial \mathbf{W}^0} = \frac{\partial L}{\partial \mathbf{y}^0} \times \frac{\partial \mathbf{y}^0}{\partial \mathbf{v}^0} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{W}^0}$$

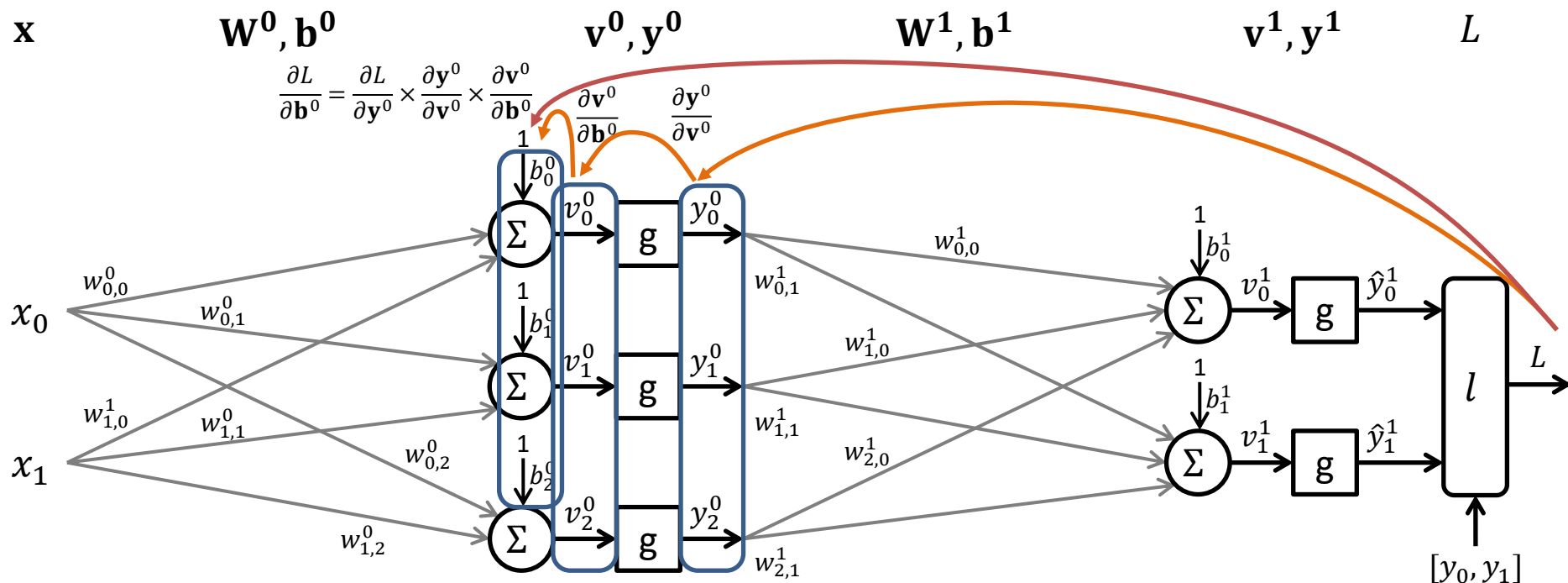
$$\frac{\partial L}{\partial \mathbf{y}^0} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{y}^0}$$

$$\frac{\partial L}{\partial \mathbf{b}^0} = \frac{\partial L}{\partial \mathbf{y}^0} \times \frac{\partial \mathbf{y}^0}{\partial \mathbf{v}^0} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{b}^0}$$

$$\frac{\partial L}{\partial \mathbf{W}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{W}^1}$$

$$\frac{\partial L}{\partial \mathbf{b}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{b}^1}$$

Backpropagation – Algorithm and equations



$$\frac{\partial L}{\partial \mathbf{W}^0} = \frac{\partial L}{\partial \mathbf{y}^0} \times \frac{\partial \mathbf{y}^0}{\partial \mathbf{v}^0} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{W}^0}$$

$$\frac{\partial L}{\partial \mathbf{y}^0} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{y}^0}$$

$$\frac{\partial L}{\partial \mathbf{b}^0} = \frac{\partial L}{\partial \mathbf{y}^0} \times \frac{\partial \mathbf{y}^0}{\partial \mathbf{v}^0} \times \frac{\partial \mathbf{v}^0}{\partial \mathbf{b}^0}$$

$$\frac{\partial L}{\partial \mathbf{W}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{W}^1}$$

$$\frac{\partial L}{\partial \mathbf{b}^1} = \frac{\partial L}{\partial \hat{\mathbf{y}}^1} \times \frac{\partial \hat{\mathbf{y}}^1}{\partial \mathbf{v}^1} \times \frac{\partial \mathbf{v}^1}{\partial \mathbf{b}^1}$$

BACKPROPAGATION – FUNCTIONS AND DERIVATIVES

Backpropagation – Functions and derivatives

- Derivative of a function:

$$- \frac{d}{dx} f(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Backpropagation – Functions and derivatives

- Sigmoid function:

- $g(v) = \frac{1}{1+e^{-v}}$

- Derivative of sigmoid function:

- $\frac{d}{dv} g(v) = g(v)(1 - g(v))$

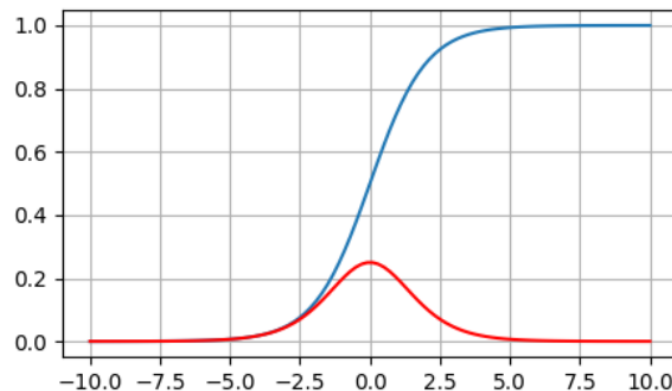
- ReLU function:

- $g(v) = \max(0, v)$ ou $g(v) = \begin{cases} v, & v > 0 \\ 0, & v \leq 0 \end{cases}$

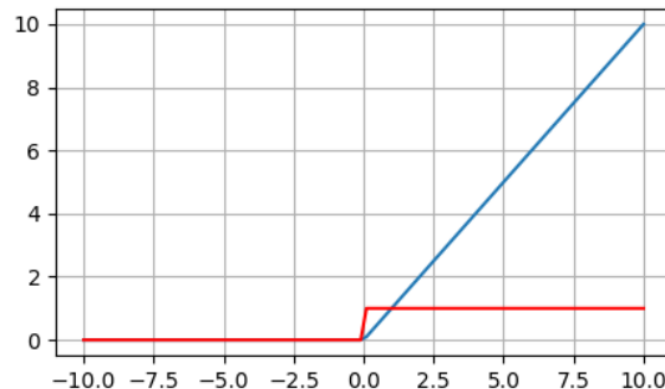
- Derivative of the ReLU function:

- $\frac{d}{dv} g(v) = \begin{cases} 1, & v > 0 \\ 0, & v \leq 0 \end{cases}$

Sigmoid function

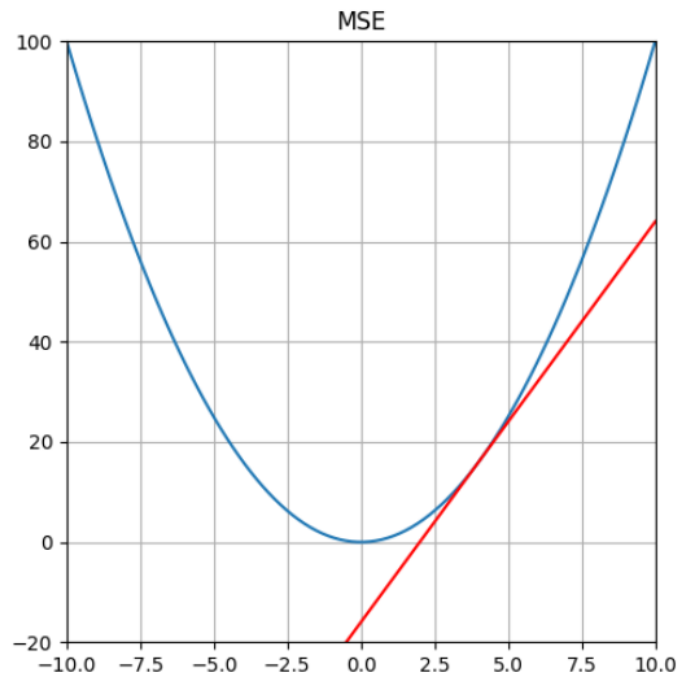


ReLU function



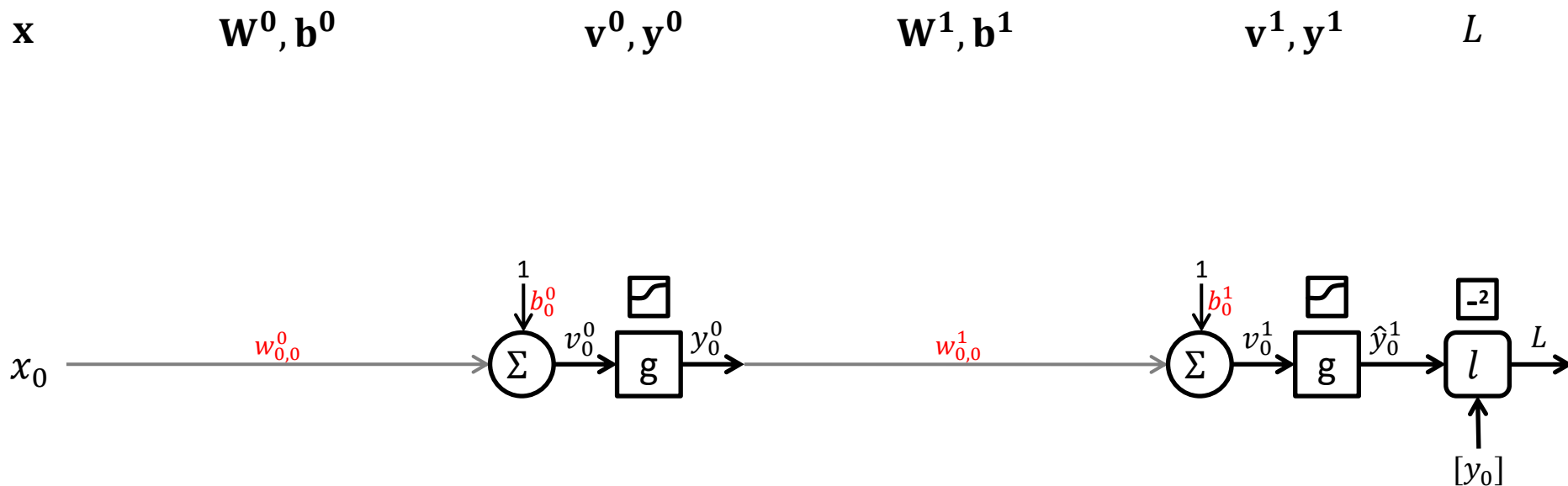
Backpropagation – Functions and derivatives

- Mean squared error (MSE):
 - $L(\hat{y}, y) = \frac{1}{N} (y - \hat{y})^2$, onde N é o número de classes.
- Partial derivative of the error function w.r.t. y :
 - $\frac{\partial L}{\partial y} = -\frac{2}{N} (y - \hat{y})$,

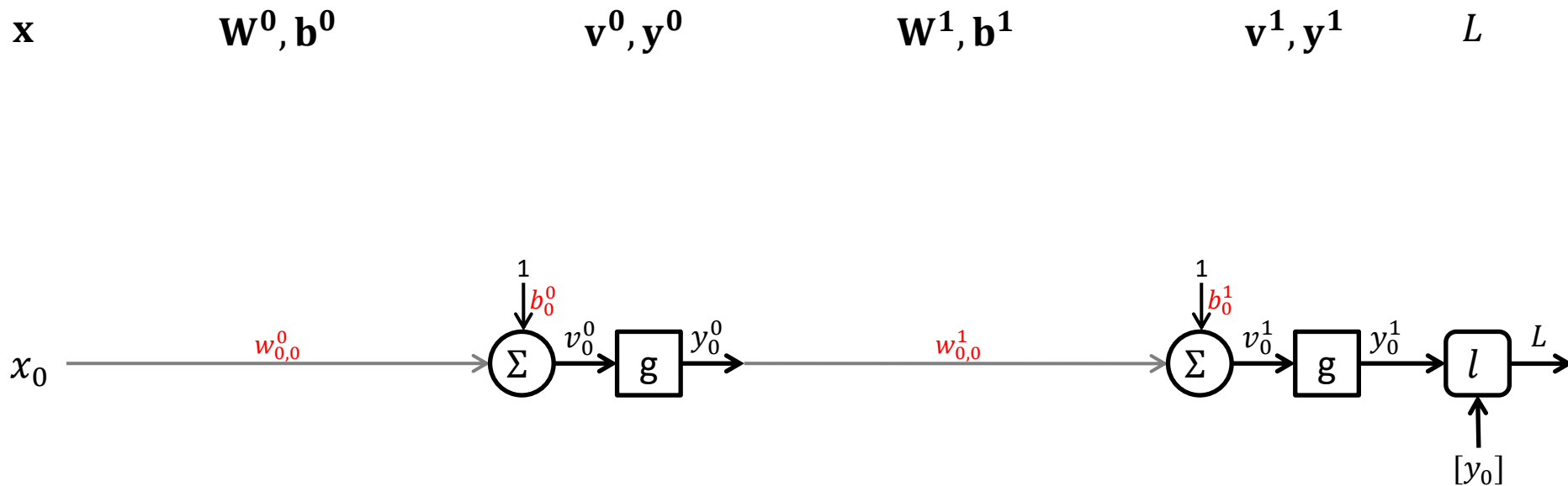


BACKPROPAGATION – EXAMPLES

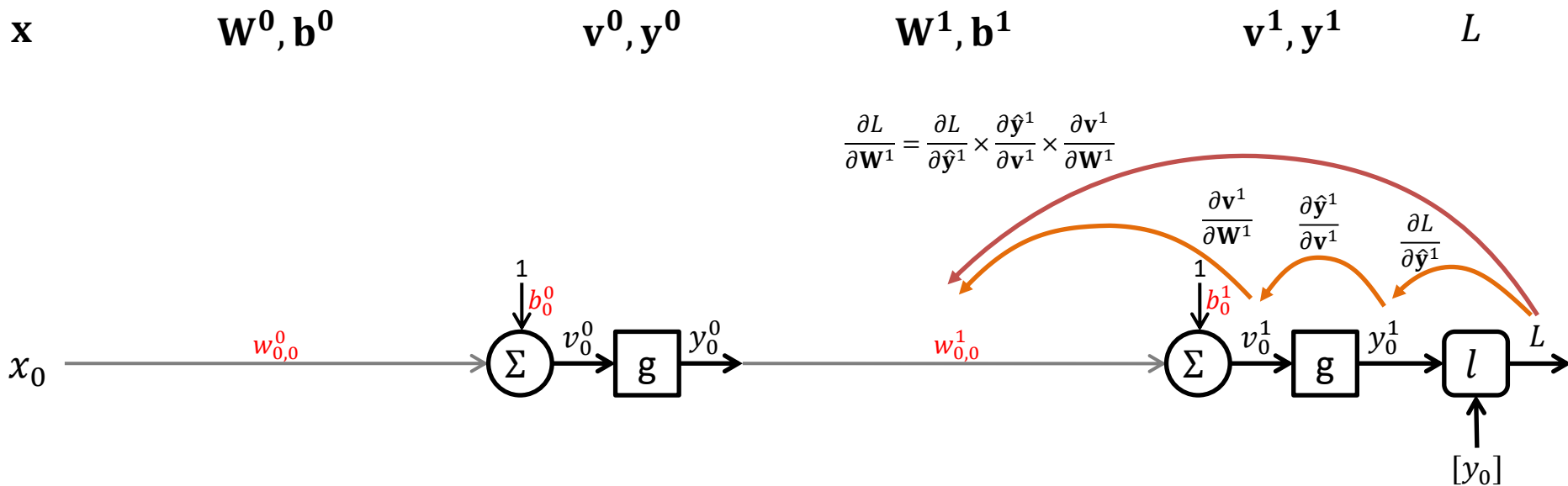
Backpropagation – Example 1



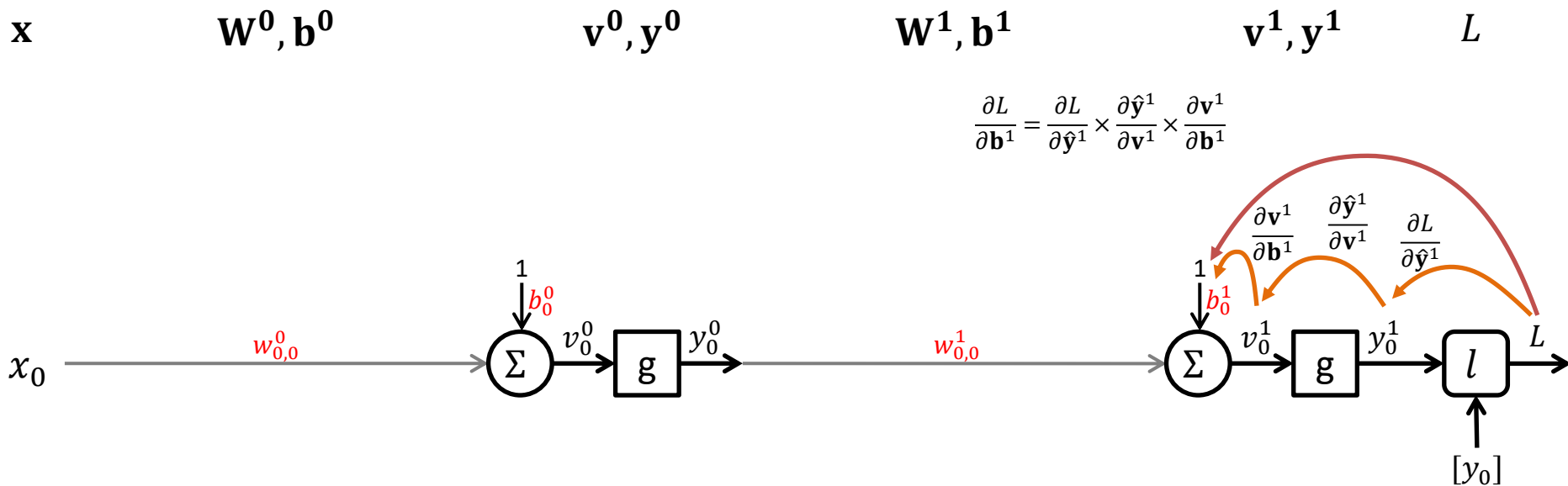
Backpropagation – Example 1



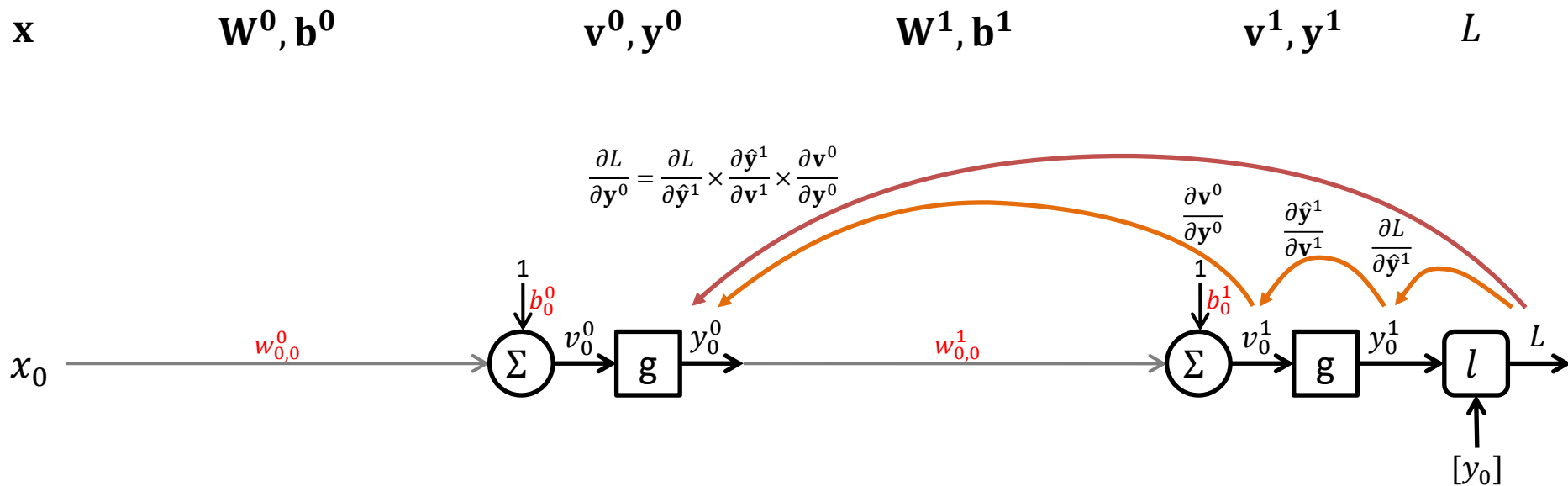
Backpropagation – Example 1



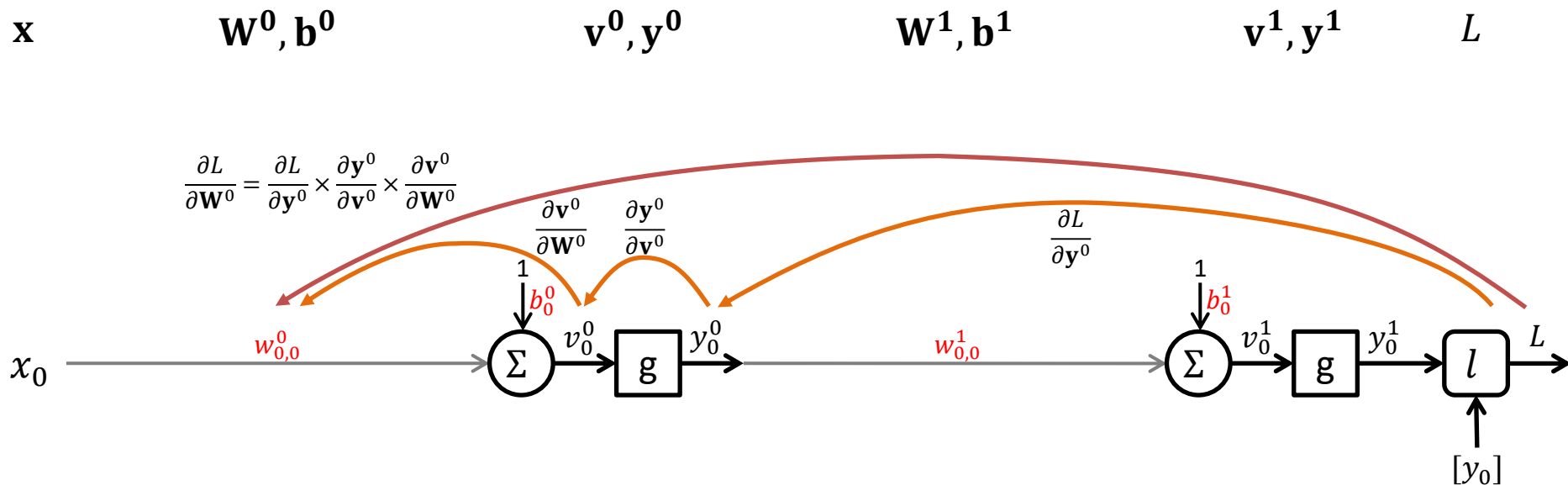
Backpropagation – Example 1



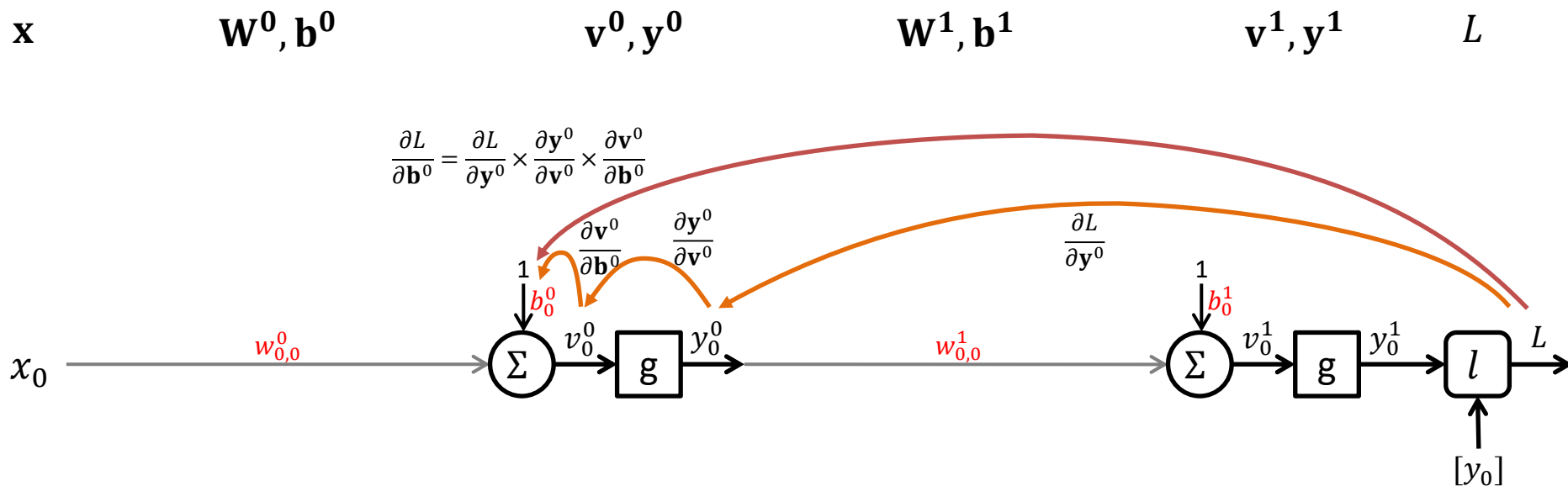
Backpropagation – Example 1



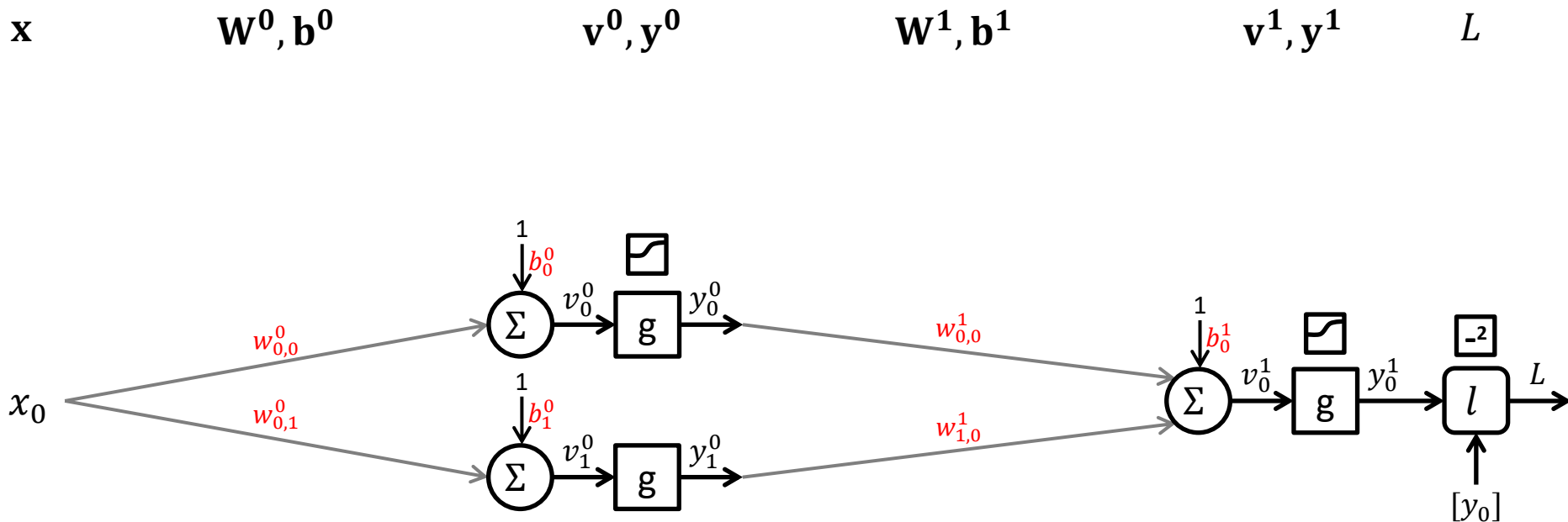
Backpropagation – Example 1



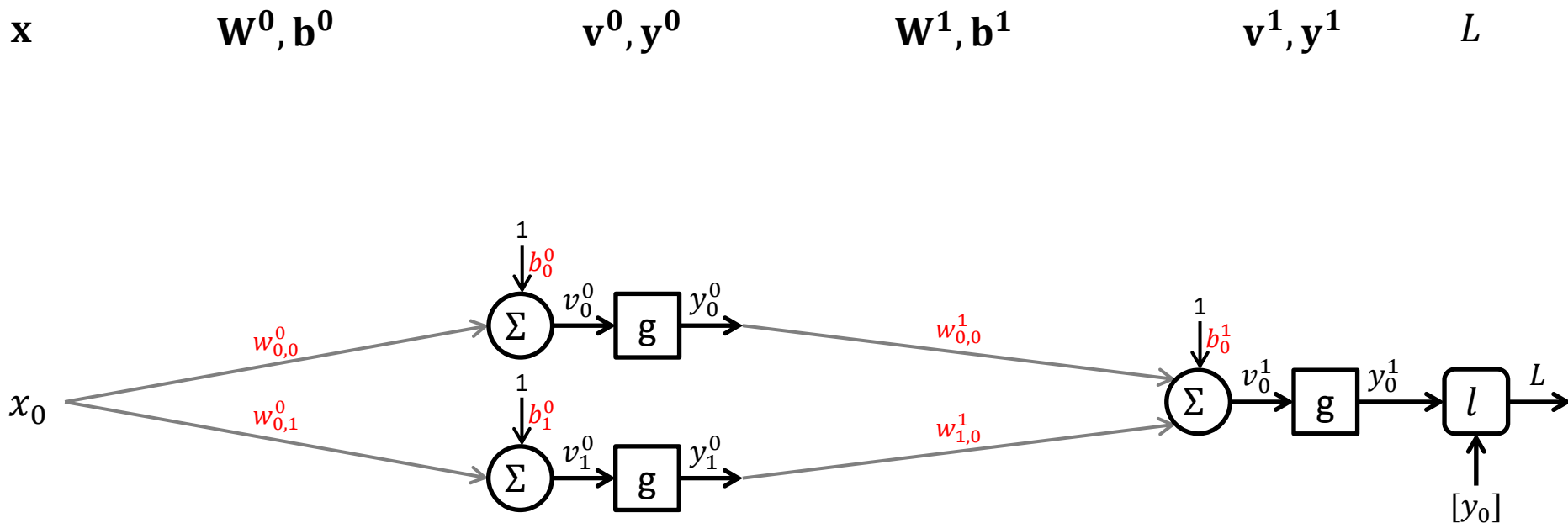
Backpropagation – Example 1



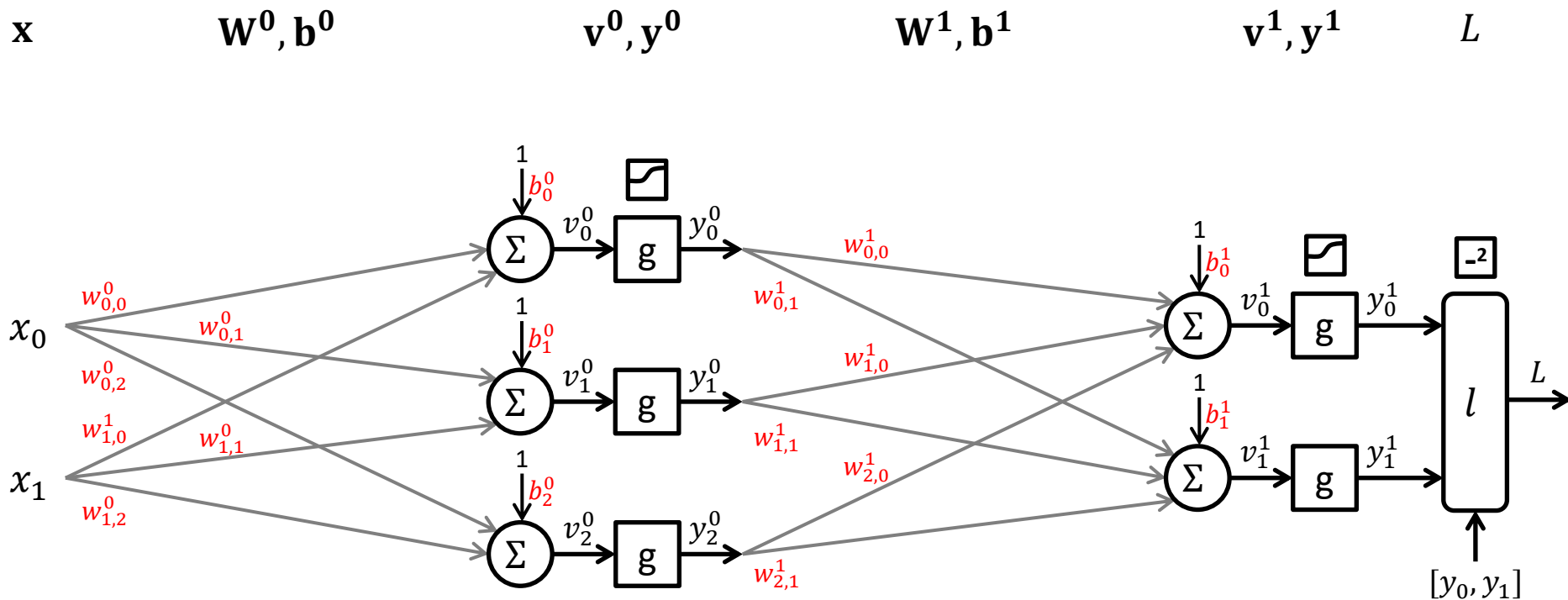
Backpropagation – Example 2



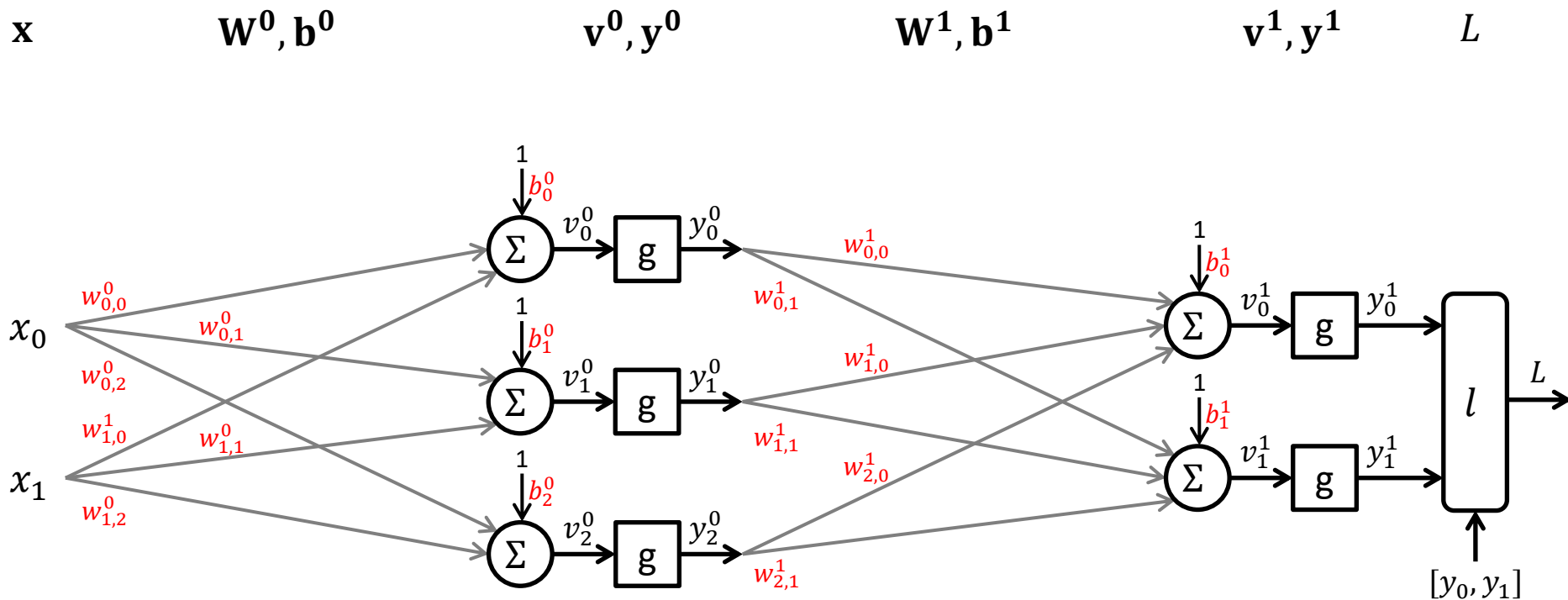
Backpropagation – Example 2



Backpropagation – Example 2



Backpropagation – Example 2



- Ponti et al. **Everything You Wanted to Know about Deep Learning for Computer Vision but Were Afraid to Ask**. Sibgrapi 2017.
 - <https://sites.icmc.usp.br/moacir/p17sibgrapi-tutorial/>
- Martin Görner. **Learn TensorFlow and deep learning, without a Ph.D.** Google Cloud.
 - <https://cloud.google.com/blog/products/gcp/learn-tensorflow-and-deep-learning-without-a-phd>
- CS231n: Convolutional Neural Networks for Visual Recognition. Stanford University.
 - <http://cs231n.stanford.edu/>
 - <http://cs231n.github.io/>
- Goodfellow, Bengio and Courville. **Deep Learning**. MIT Press, 2016
 - <https://www.deeplearningbook.org/>

- Rabindra Lamsal. A step by step forward pass and backpropagation example
 - <https://theneuralblog.com/forward-pass-backpropagation-example/>
- Matt Mazur. A Step by Step Backpropagation Example
 - <https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>

Bibliography

- Back-Propagation is very simple. Who made it Complicated ?
 - <https://medium.com/@14prakash/back-propagation-is-very-simple-who-made-it-complicated-97b794c97e5c>
- Chapter 7: Artificial neural networks with Math.
 - <https://medium.com/deep-math-machine-learning-ai/chapter-7-artificial-neural-networks-with-math-bb711169481b>
- The Matrix Calculus You Need For Deep Learning
 - <http://explained.ai/matrix-calculus/index.html>
- How backpropagation works, and how you can use Python to build a neural network
 - <https://medium.freecodecamp.org/build-a-flexible-neural-network-with-backpropagation-in-python-acffeb7846d0>
- All the Backpropagation derivatives
 - <https://medium.com/@pdquant/all-the-backpropagation-derivatives-d5275f727f60>
- Brent Scarff. **Understanding Backpropagation.**
 - <https://towardsdatascience.com/understanding-backpropagation-abcc509ca9d0>

APPENDIX: EXAMPLE

Example 1

FORWARD:

$$v_0 = x_0 w_{00} + b_0 = 0.7 \times 0.1 + 0.25 = 0.32$$

$$y_0 = \frac{1}{1 + 0.7261} = 0.5793$$

$$MSE = L = (1 - 0.6547)^2 = 0.1192$$

Sigmoide:

$$g(v) = \frac{1}{1 + e^{-v}}$$

Derivada:

$$\frac{d}{dv} g(v) = g(v)(1 - g(v))$$

MSE

$$L(y, \hat{y}) = \frac{1}{N} (y - \hat{y})^2$$

Derivada parcial:

$$\frac{\partial L}{\partial y} = -\frac{2}{N} (y - \hat{y})$$

FORWARD:

$$\frac{\partial L}{\partial w_{00}} = \frac{\partial L}{\partial y_0} \times \frac{\partial y_0}{\partial v_0} \times \frac{\partial v_0}{\partial w_{00}} = -0.6906 \times 0.2261 \times 0.1 = -0.0156$$

$$\frac{\partial L}{\partial y_0} = -2(1 - 0.6547) = -0.6906$$

$$\frac{\partial y_0}{\partial v_0} = 0.6547 \times (1 - 0.6547) = 0.2261$$

$$\frac{\partial v_0}{\partial w_{00}} = 0.7$$

$$\frac{\partial L}{\partial w_{01}} = -0.6906 \times 0.2261 \times 0.5 = -0.0778$$

$$\frac{\partial L}{\partial b_0} = -0.6906 \times 0.2261 \times 1 = -0.1561$$

Sigmoide:

$$g(v) = \frac{1}{1 + e^{-v}}$$

Derivada:

$$\frac{d}{dv} g(v) = g(v)(1 - g(v))$$

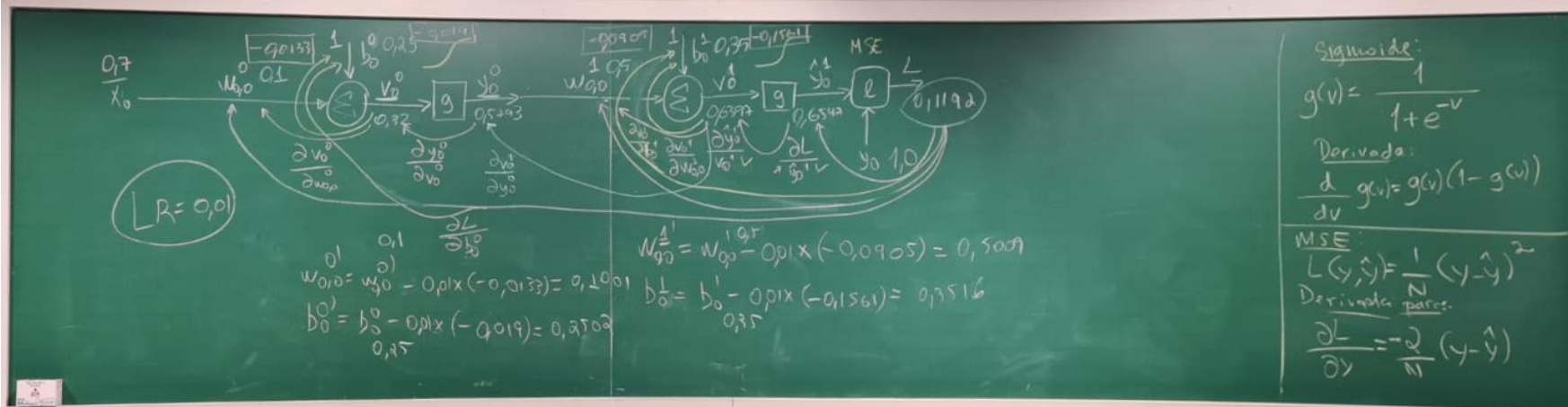
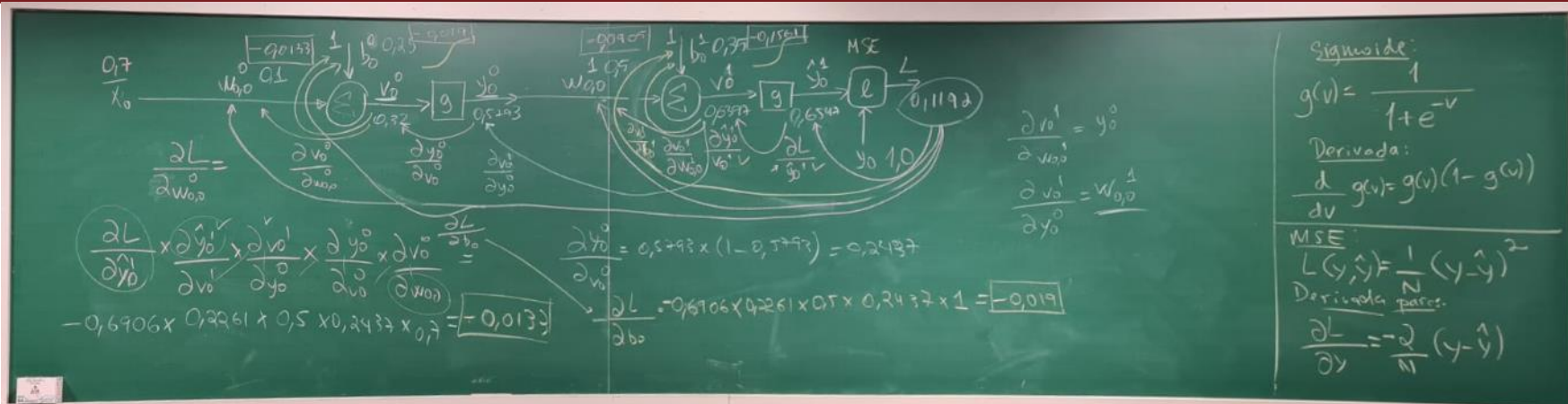
MSE

$$L(y, \hat{y}) = \frac{1}{N} (y - \hat{y})^2$$

Derivada parcial:

$$\frac{\partial L}{\partial y} = -\frac{2}{N} (y - \hat{y})$$

Example 1



THE END