

COMP90051

Workshop Week 05

About the Workshops

- 7 sessions in total
 - Tue 12:00-13:00 AH211
 - Tue 12:00-13:00 AH108 *
 - Tue 13:00-14:00 AH210
 - Tue 16:15-17:15 AH109
 - Tue 17:15-18:15 AH236 *
 - Tue 18:15-19:15 AH236 *
 - Fri 14:15-15:15 AH211

About the Workshops

- Homepage

- <https://trevorcohn.github.io/comp90051-2017/workshops>

- Solutions will be released on next Friday (a week later).

Syllabus

1	Introduction; Probability theory	Probabilistic models; Parameter fitting	
2	Linear regression; Intro to regularization	Logistic regression; Basis expansion	
3	Optimization; Regularization	Perceptron	
4	Backpropagation	CNNs; Auto-encoders	←
5	Hard-margin SVMs	Soft-margin SVMs	
6	Additional topics	Kernel methods	
7	Unsupervised learning	Unsupervised learning	
8	Dimensionality reduction; Principal component analysis	Multidimensional scaling; Spectral clustering	
9	Bayesian fundamentals	Bayesian inference with conjugate priors	
10	PGMs, fundamentals	Conditional independence	
11	PGMs, inference	Belief propagation	
12	Statistical inference; Apps	Subject review	

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Gradient descent & stochastic gradient descent (SGD)
 - ❑ Gradient and backpropagation
 - ❑ Logistic regression
 - ❑ Neural networks with one hidden layer
- ❑ Notebook tasks
 - ❑ Task 1: Multi-layer perceptron, SGD

Outline

- Review the lecture, background knowledge, etc.
 - Gradient descent & stochastic gradient descent (SGD)
 - Gradient and backpropagation
 - Logistic regression
 - Neural networks with one hidden layer
- Notebook tasks
 - Task 1: Multi-layer perceptron, SGD

Gradient descent & Stochastic GD (SGD)

- To minimize an objective function $obj(w)$
- Usually, obj is the average loss plus a regularization term

$$\min_{\mathbf{w}} obj(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(f(x_i; \mathbf{w}), y_i) + \lambda R(\mathbf{w})$$

Gradient descent

- To minimize an objective function $obj(w)$
- Usually, obj is the average loss plus a regularization term

$$\min_{\mathbf{w}} obj(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(f(x_i; \mathbf{w}), y_i) + \lambda R(\mathbf{w})$$

- Loop until \mathbf{w} doesn't change

$$\mathbf{w} = \mathbf{w} - \eta \frac{\partial obj(\mathbf{w})}{\partial \mathbf{w}}$$

Gradient descent

- To minimize an objective function $obj(w)$
- Usually, obj is the average loss plus a regularization term

$$\min_{\mathbf{w}} obj(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(f(x_i; \mathbf{w}), y_i) + \lambda R(\mathbf{w})$$

- Loop until \mathbf{w} doesn't change

$$\text{grad}_{\mathbf{w}} = \frac{\partial obj(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L(f(x_i; \mathbf{w}), y_i)}{\partial \mathbf{w}} + \lambda \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$$
$$\mathbf{w} = \mathbf{w} - \eta \text{grad}_{\mathbf{w}}$$

Stochastic gradient descent (SGD)

- To minimize an objective function $obj(w)$
- Usually, obj is the average loss plus a regularization term

$$\min_{\mathbf{w}} obj(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(f(x_i; \mathbf{w}), y_i) + \lambda R(\mathbf{w})$$

- Loop until \mathbf{w} doesn't change

- Sample i from $\{1, 2, \dots, N\}$ or For $i = 1, 2, \dots, N$

$$\text{grad}_{\mathbf{w}} = \frac{\partial L(f(x_i; \mathbf{w}), y_i)}{\partial \mathbf{w}} + \lambda \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$$
$$\mathbf{w} = \mathbf{w} - \eta \text{grad}_{\mathbf{w}}$$

Stochastic gradient descent (SGD)

□ Note: SGD has other variants

□ Loop until \mathbf{w} doesn't change ← online learning

□ Sample i from $\{1, 2, \dots, N\}$ or For $i = 1, 2, \dots, N$

$$\text{grad}_{\mathbf{w}} = \frac{\partial L(f(x_i; \mathbf{w}), y_i)}{\partial \mathbf{w}} + \lambda \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$$
$$\mathbf{w} = \mathbf{w} - \eta \text{grad}_{\mathbf{w}}$$

□ Loop until \mathbf{w} doesn't change ← mini-batch

□ Sample a subset S from $\{1, 2, \dots, N\}$

$$\text{grad}_{\mathbf{w}} = \frac{1}{|S|} \sum_{i \in S} \frac{\partial L(f(x_i; \mathbf{w}), y_i)}{\partial \mathbf{w}} + \lambda \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$$
$$\mathbf{w} = \mathbf{w} - \eta \text{grad}_{\mathbf{w}}$$

Gradient descent & Stochastic GD (SGD)

□ Gradient descent

□ Loop until \mathbf{w} doesn't change

$$\text{grad}_{\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial L(f(x_i; \mathbf{w}), y_i)}{\partial \mathbf{w}} + \lambda \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$$

$$\mathbf{w} = \mathbf{w} - \eta \text{grad}_{\mathbf{w}}$$

□ Stochastic GD (SGD)

□ Loop until \mathbf{w} doesn't change

□ Sample i from $\{1, 2, \dots, N\}$
or For $i = 1, 2, \dots, N$

$$\text{grad}_{\mathbf{w}} = \frac{\partial L(f(x_i; \mathbf{w}), y_i)}{\partial \mathbf{w}} + \lambda \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$$

$$\mathbf{w} = \mathbf{w} - \eta \text{grad}_{\mathbf{w}}$$

Outline

- Review the lecture, background knowledge, etc.
 - Gradient descent & stochastic gradient descent (SGD)
 - Gradient and backpropagation
 - Logistic regression
 - Neural networks with one hidden layer
- Notebook tasks
 - Task 1: Multi-layer perceptron, SGD

Formulas you need to know

□ Logistic function

$$y = \sigma(x) = \frac{1}{1 + e^{-x}}$$
$$\frac{dy}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = y(1 - y)$$

□ Hyperbolic tangent function

$$y = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$
$$\frac{dy}{dx} = \frac{4}{(e^x + e^{-x})^2} = 1 - y^2$$

Formulas you need to know

□ Log-loss

$$L(\mathbf{x}_i, y_i; \mathbf{W}) = -\log p(y = y_i | \mathbf{x} = \mathbf{x}_i; \mathbf{W})$$

□ Log-loss for binary classification

$$\hat{y}_i = p(y = 1 | \mathbf{x} = \mathbf{x}_i; \mathbf{W})$$

$$L(\mathbf{x}_i, y_i; \mathbf{W}) = -(1 - y_i) \log(1 - \hat{y}_i) - y_i \log \hat{y}_i$$

$$L(\mathbf{x}_i, y_i; \mathbf{W}) = \begin{cases} -\log(1 - \hat{y}_i) & y_i = 0 \\ -\log \hat{y}_i & y_i = 1 \end{cases}$$

Formulas you need to know

□ Logistic regression (2-D points, 2 classes)

□ $\mathbf{x} = [x_1 \quad x_2] \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

□ Decision function

$$s = f(\mathbf{x}; \mathbf{w}, b) = x_1 w_1 + x_2 w_2 + b$$

□ Probability output

$$\hat{y} = \sigma(s) = \frac{1}{1 + e^{-s}}$$

□ Log-loss

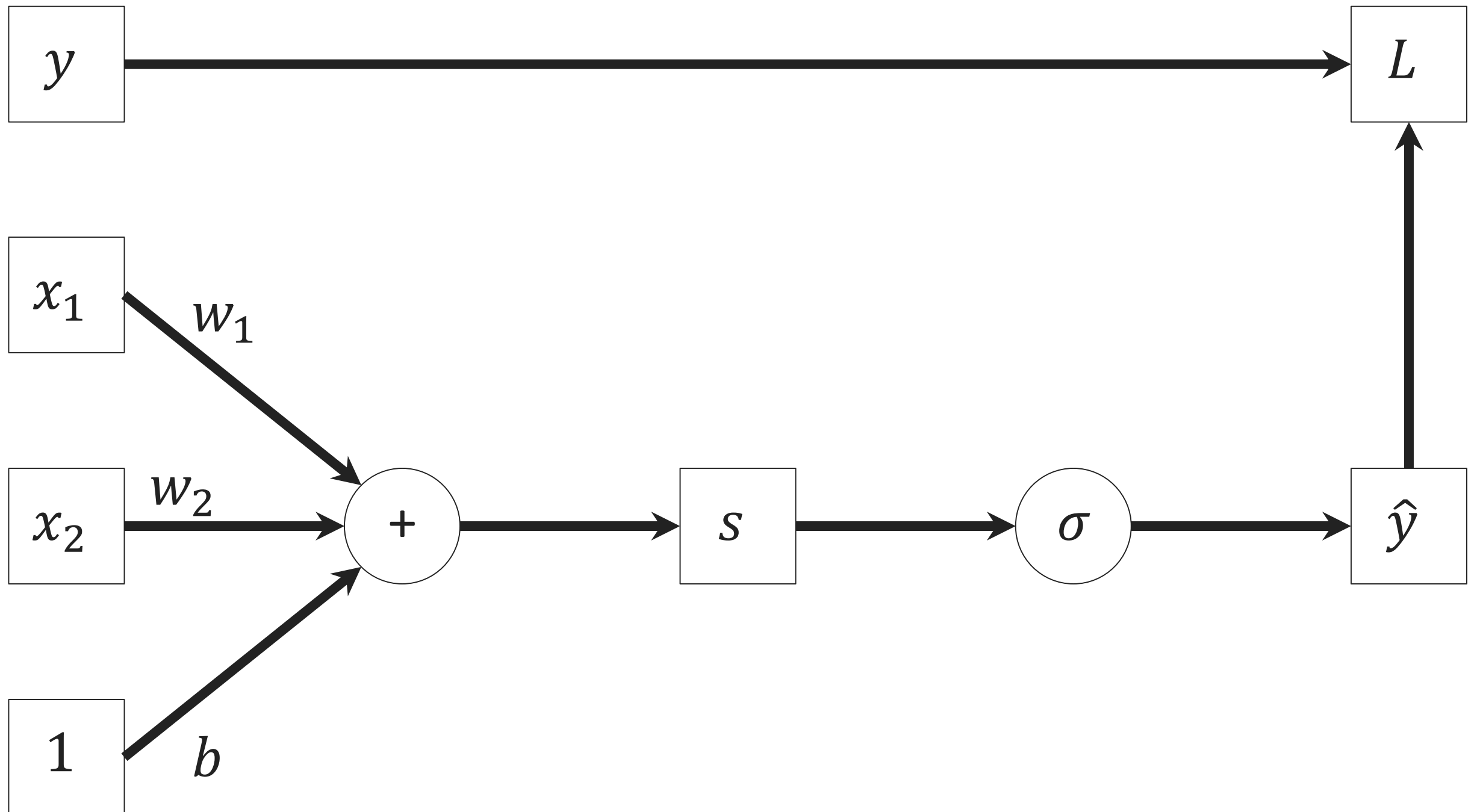
$$L(\mathbf{x}, y; \mathbf{w}, b) = -(1 - y) \log(1 - \hat{y}) - y \log \hat{y}$$

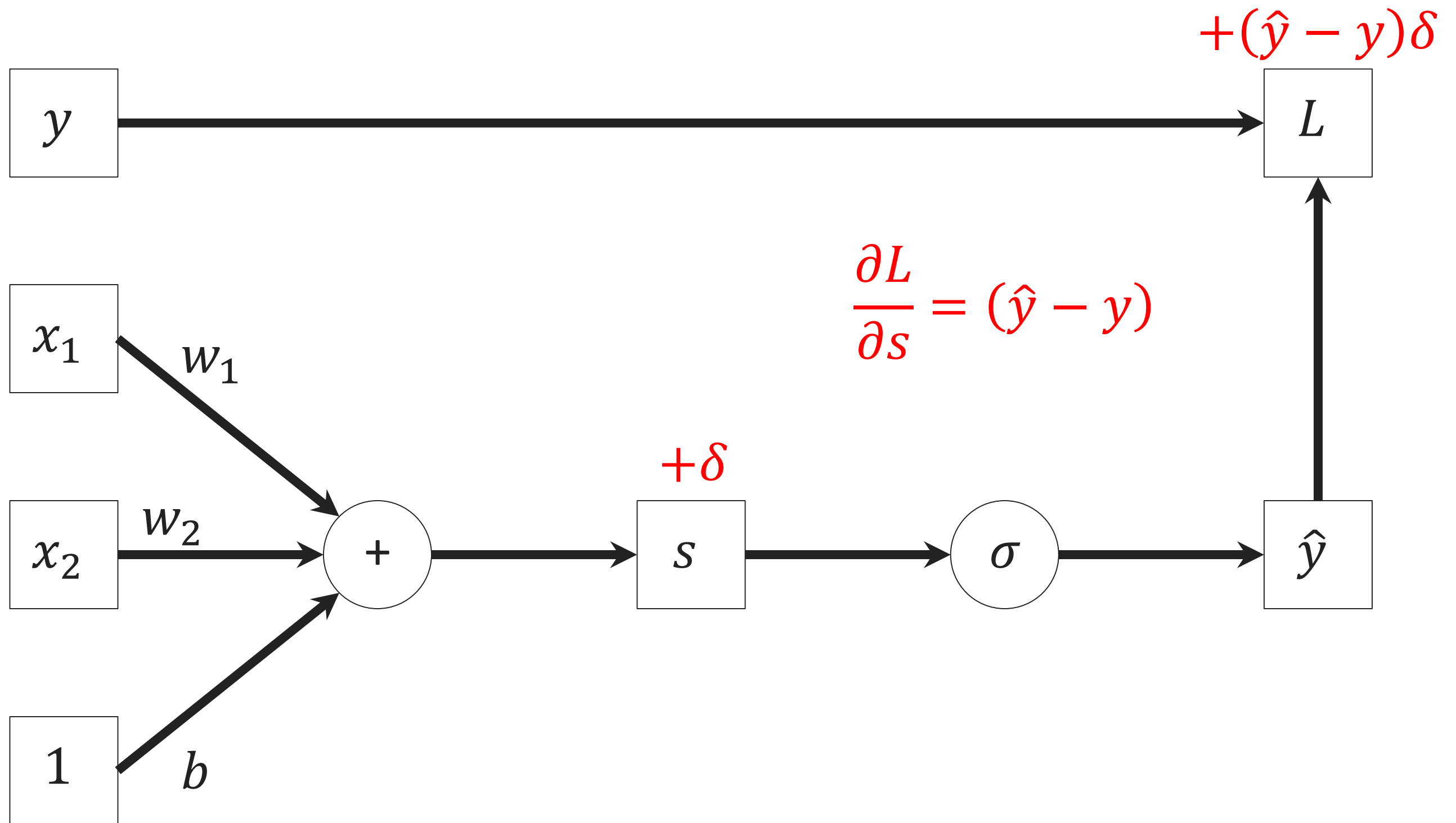
$$\frac{\partial L}{\partial s} = \frac{1}{1 + e^{-s}} - y = \hat{y} - y$$

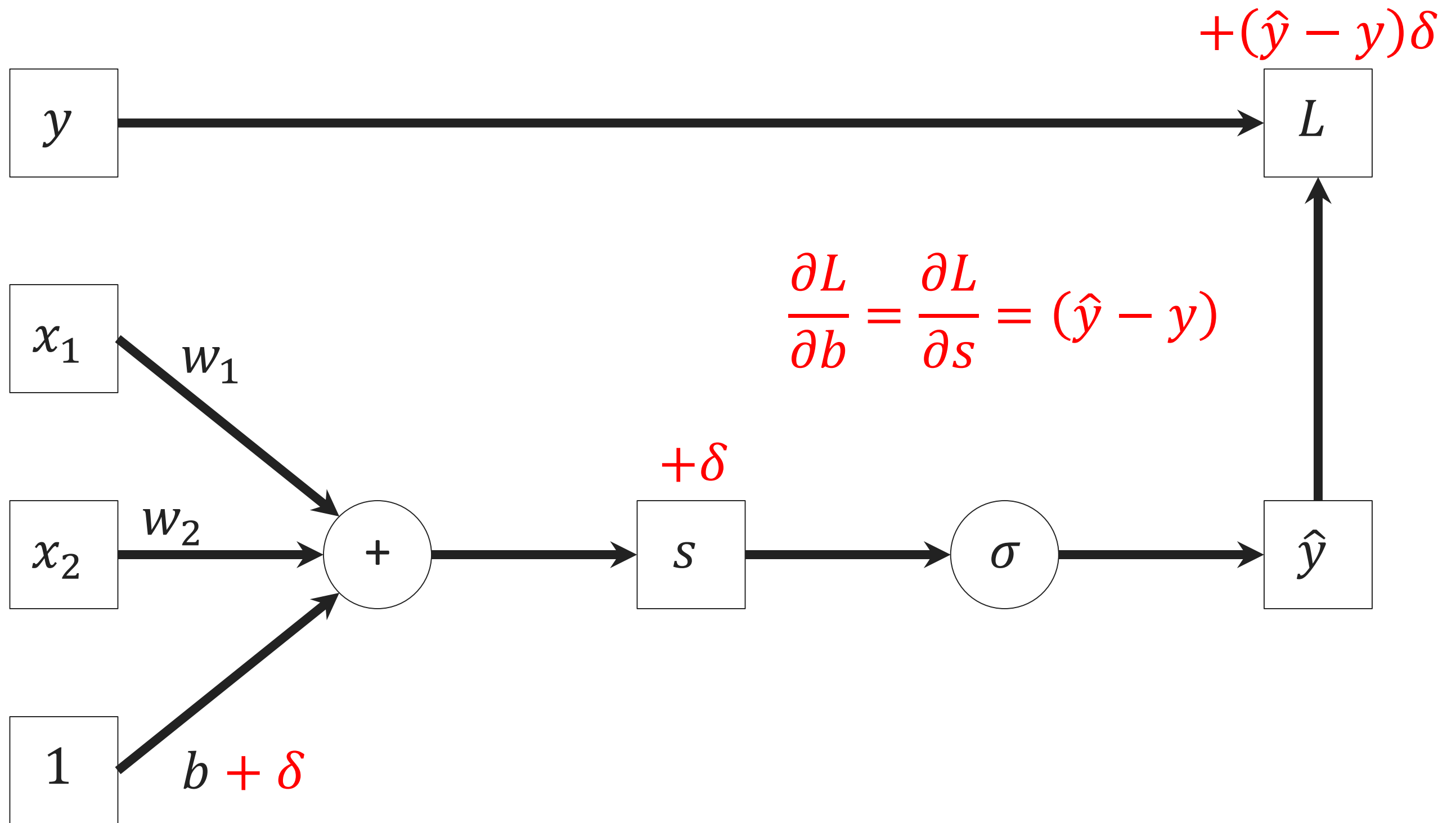
Outline

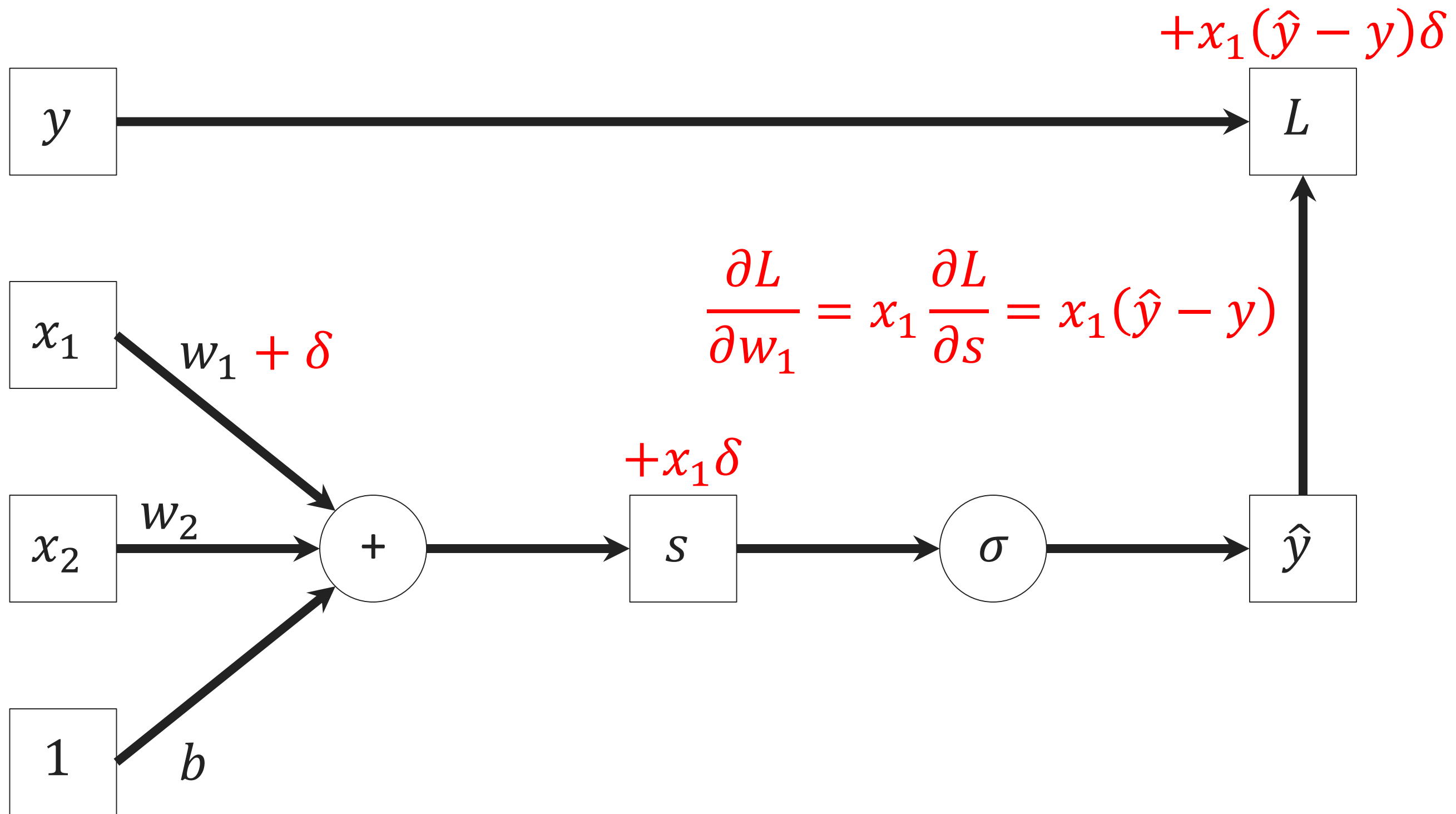
- Review the lecture, background knowledge, etc.
 - Gradient descent & stochastic gradient descent (SGD)
 - Gradient and backpropagation
 - Logistic regression
 - Neural networks with one hidden layer
- Notebook tasks
 - Task 1: Multi-layer perceptron, SGD

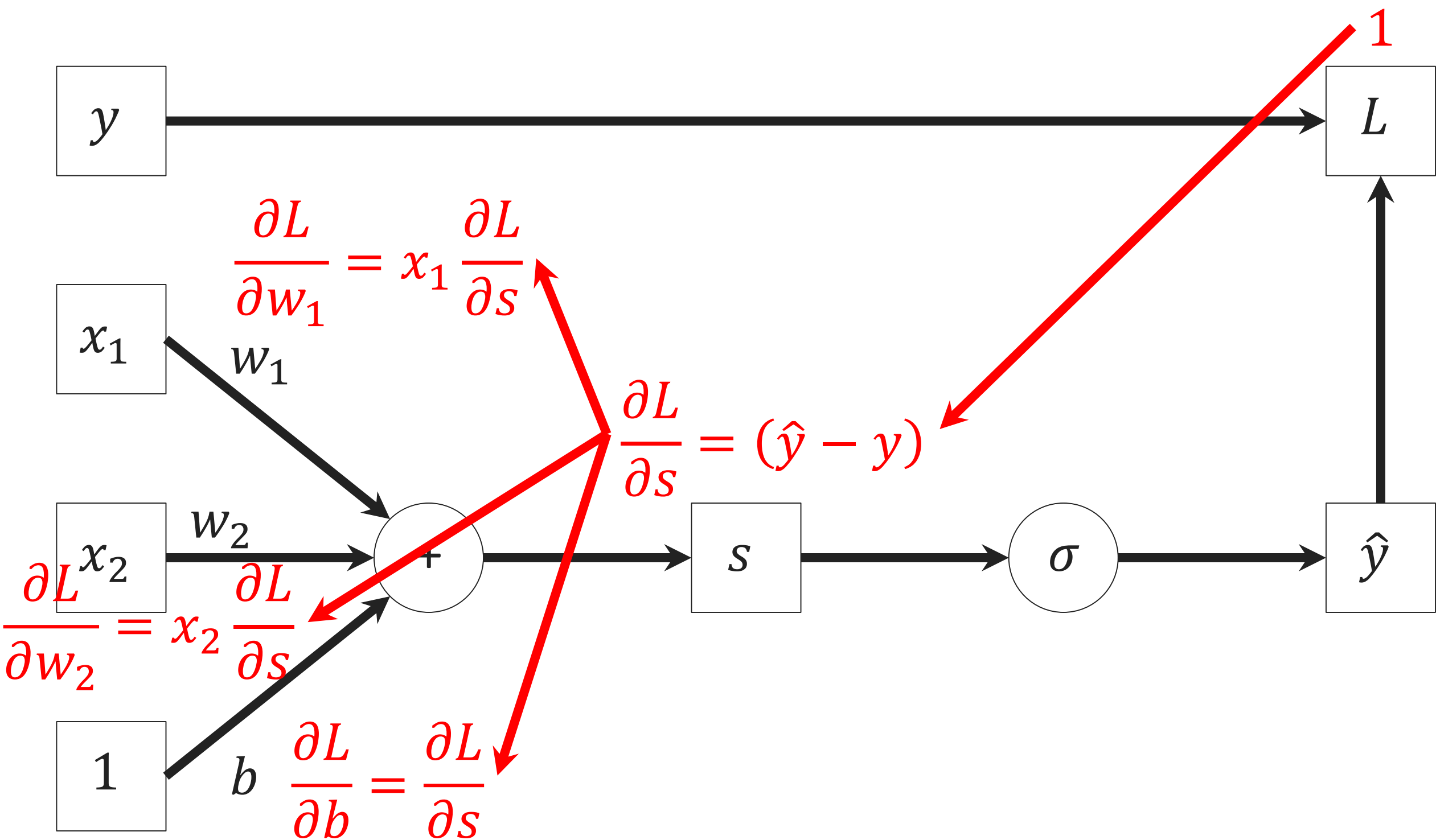
Forward pass

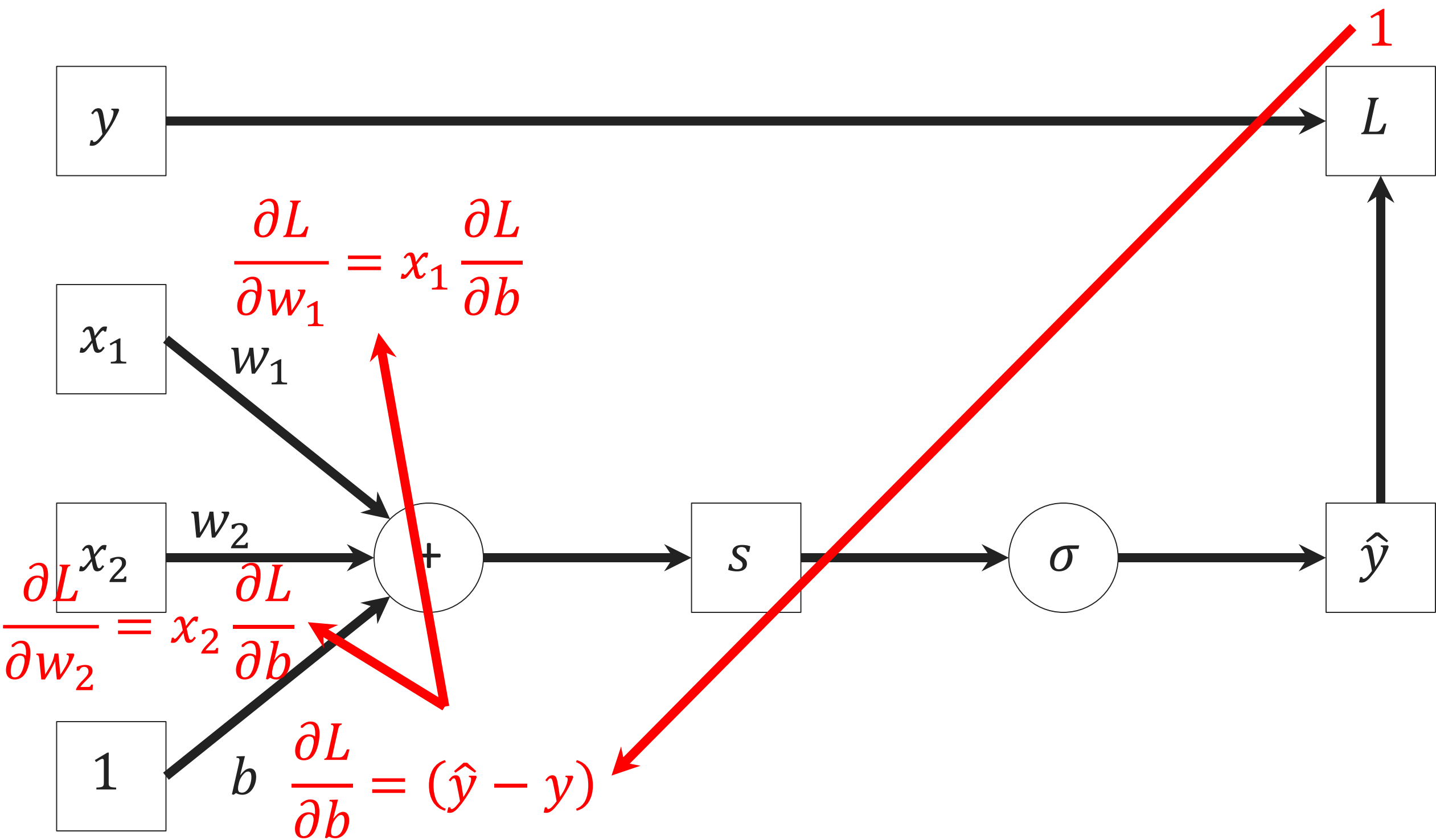




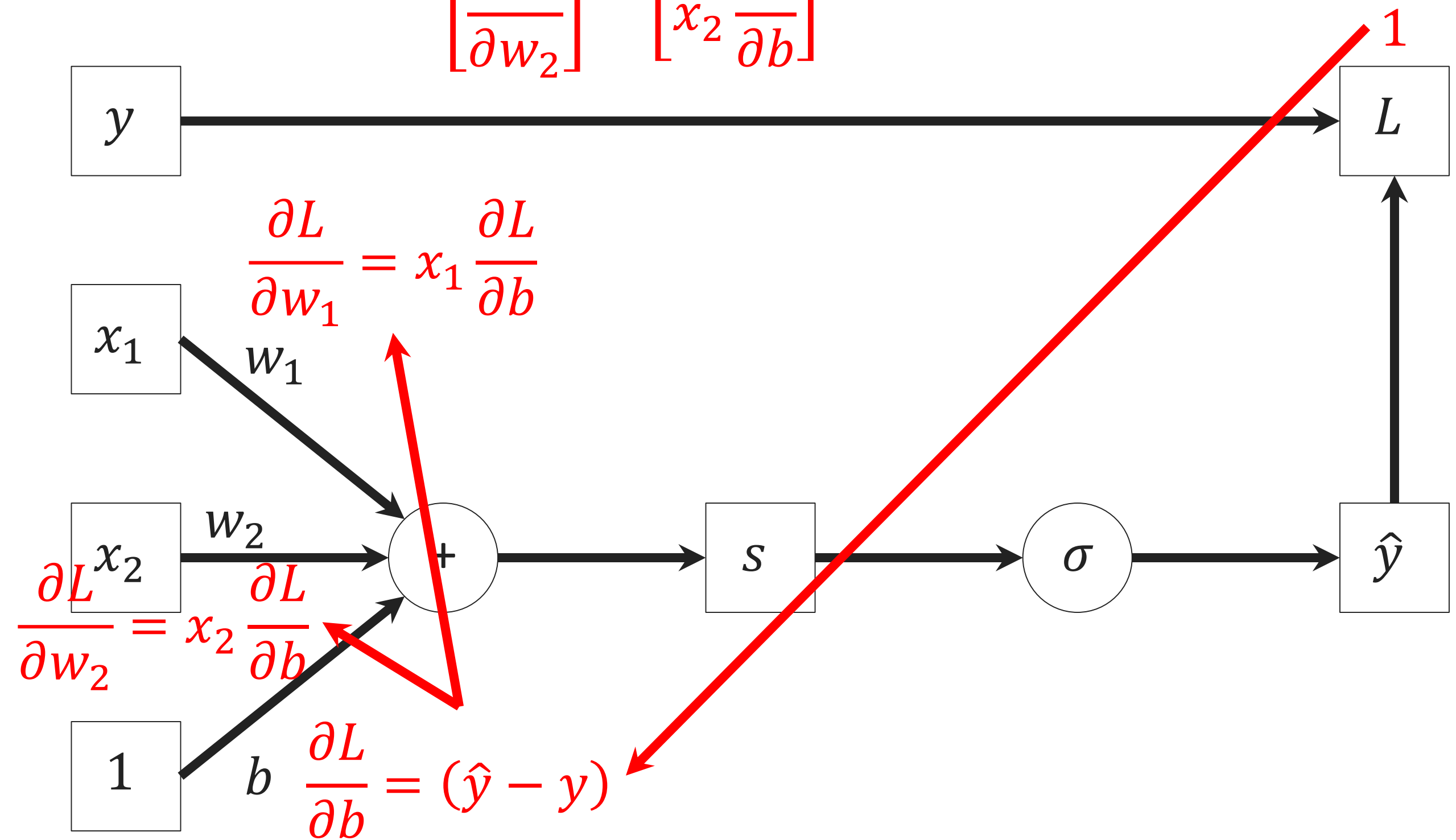




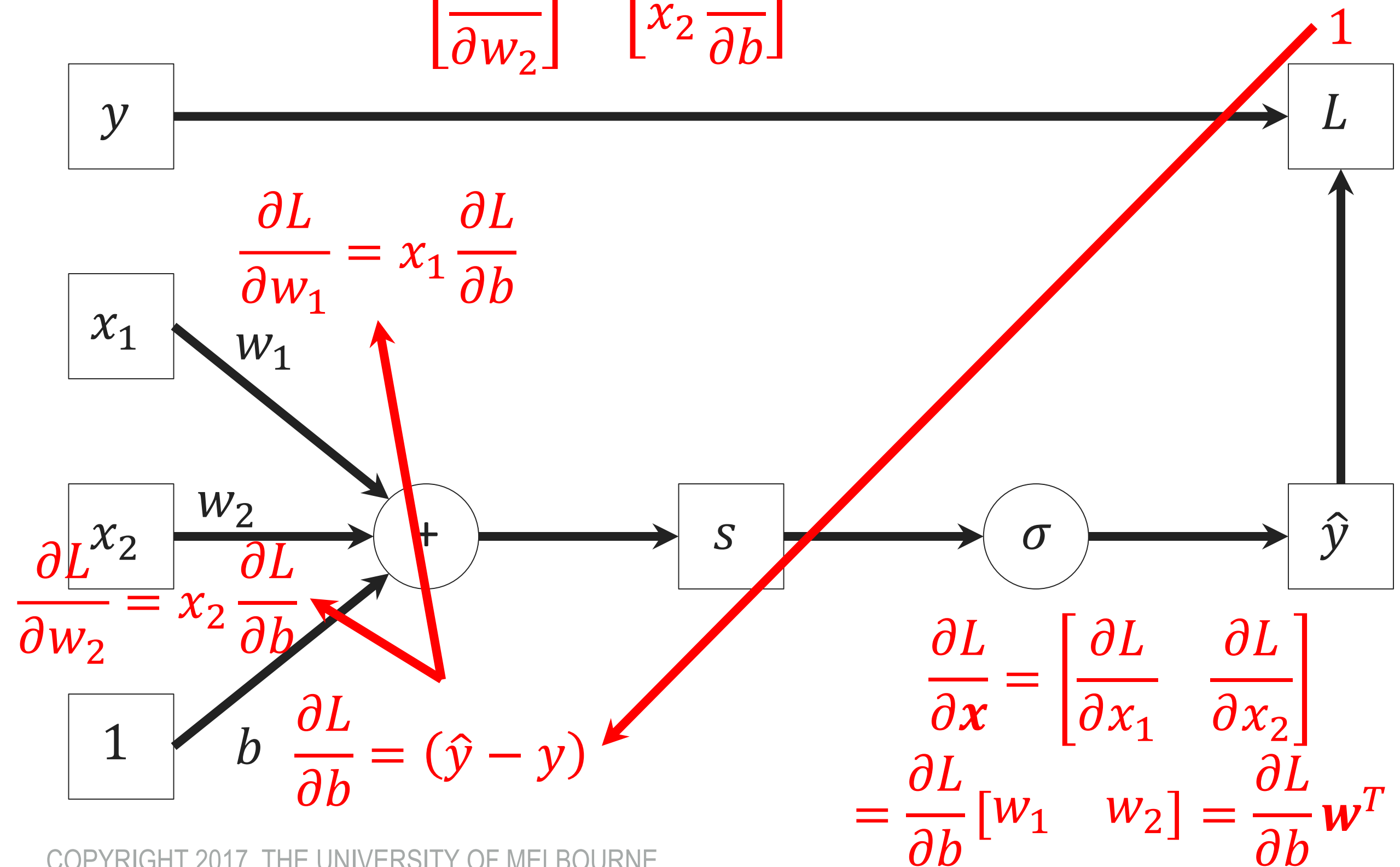




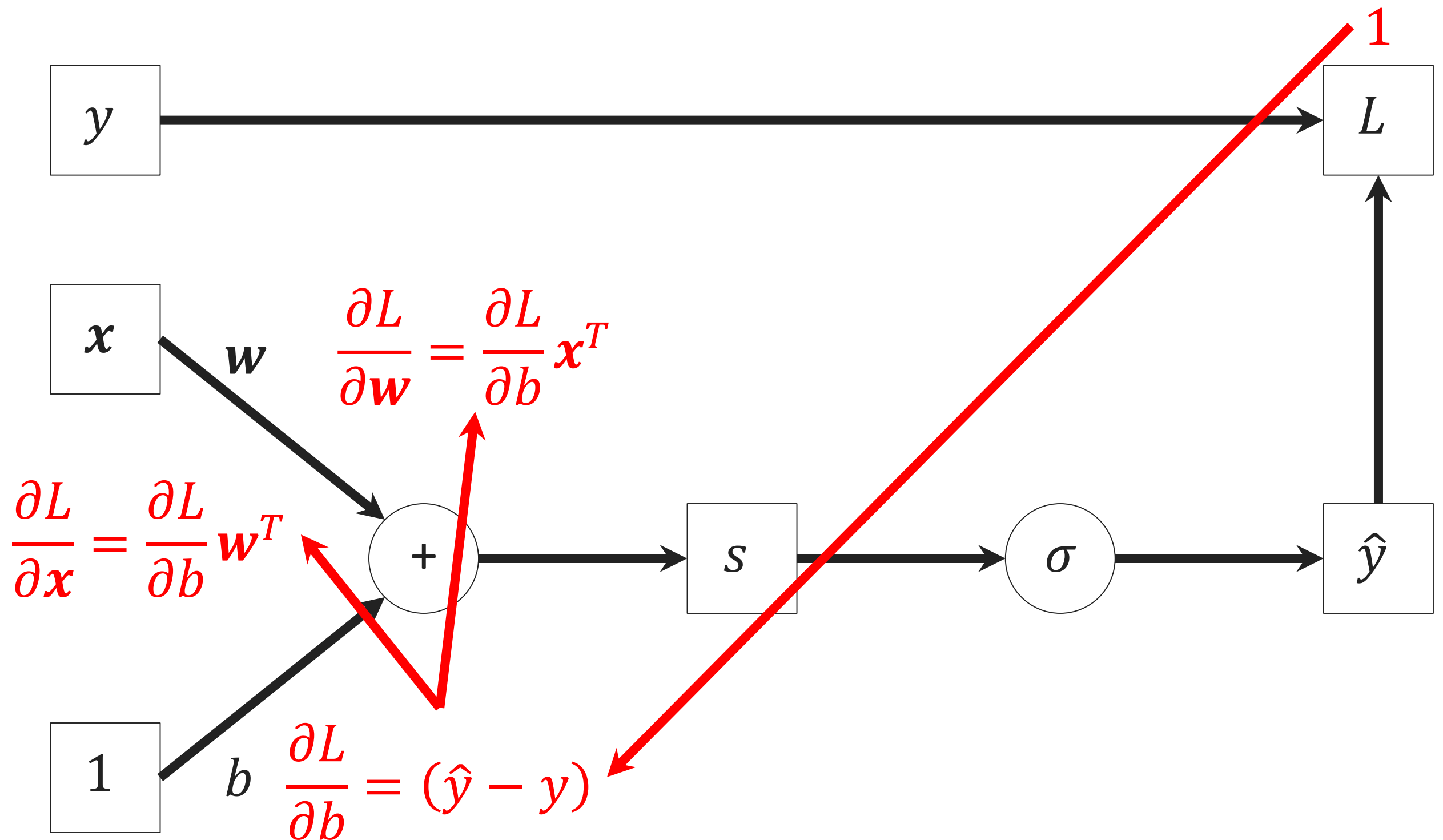
$$\frac{\partial L}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = \begin{bmatrix} x_1 \frac{\partial L}{\partial b} \\ x_2 \frac{\partial L}{\partial b} \end{bmatrix} = \frac{\partial L}{\partial b} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \frac{\partial L}{\partial b} \mathbf{x}^T$$



$$\frac{\partial L}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = \begin{bmatrix} x_1 \frac{\partial L}{\partial b} \\ x_2 \frac{\partial L}{\partial b} \end{bmatrix} = \frac{\partial L}{\partial b} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \frac{\partial L}{\partial b} \mathbf{x}^T$$



Backpropagation



Outline

- Review the lecture, background knowledge, etc.
 - Gradient descent & stochastic gradient descent (SGD)
 - Gradient and backpropagation
 - Logistic regression
 - Neural networks with one hidden layer
- Notebook tasks
 - Task 1: Multi-layer perceptron, SGD

Neural networks with one hidden layer

□ Input: 2-D points $\mathbf{x} = [x_1 \quad x_2]$

□ Hidden layer: 2 units $\mathbf{h} = [h_1 \quad h_2]$

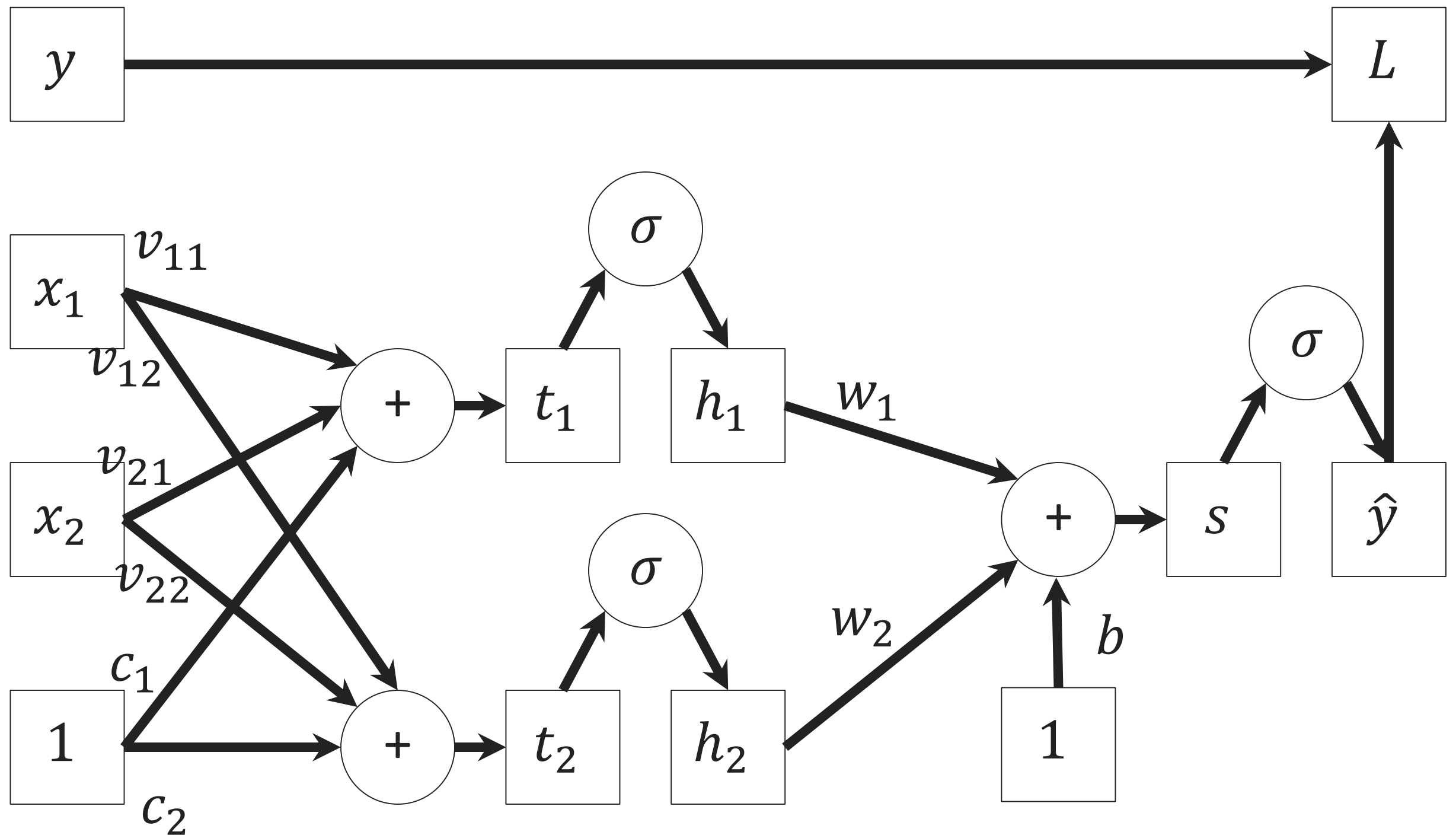
$$\mathbf{t} = [t_1 \quad t_2] = \mathbf{xV} + \mathbf{c} = [x_1 \quad x_2] \begin{bmatrix} v_{11} & v_{12} \\ v_{21} & v_{22} \end{bmatrix} + [c_1 \quad c_2]$$

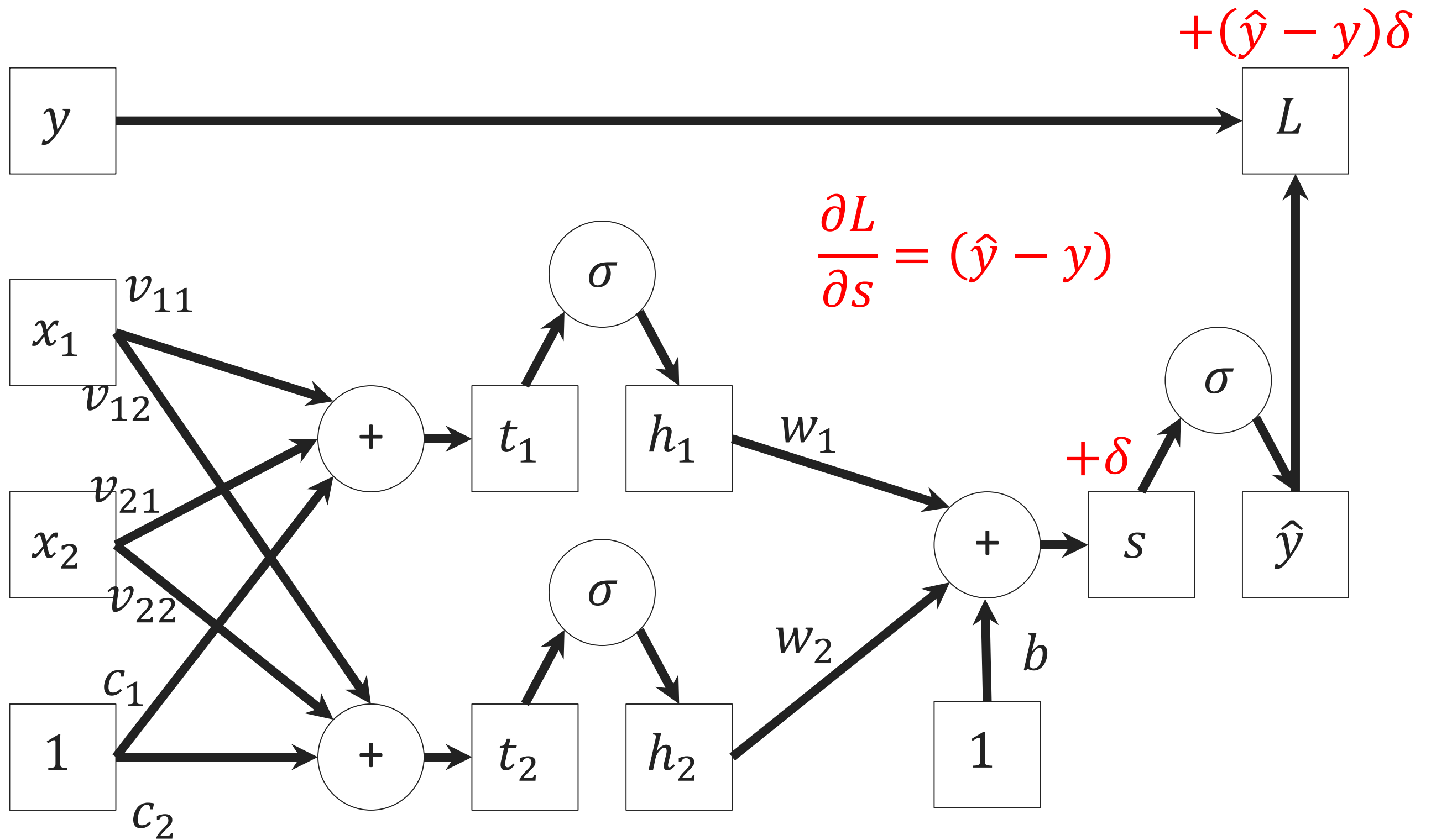
$$\mathbf{h} = [h_1 \quad h_2] = \sigma(\mathbf{t}) = \sigma([t_1 \quad t_2]) = [\sigma(t_1) \quad \sigma(t_2)]$$

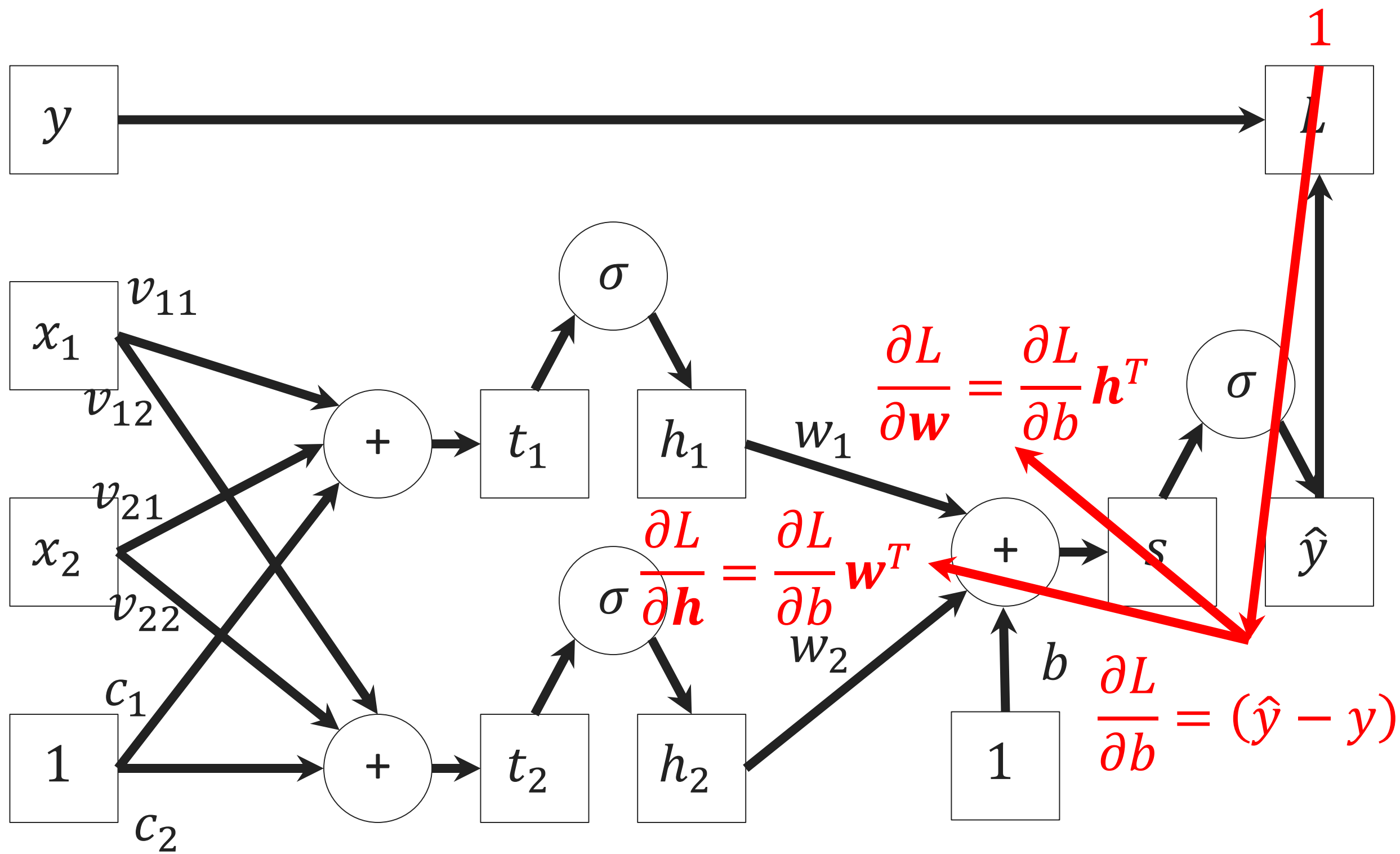
□ Output: \hat{y}

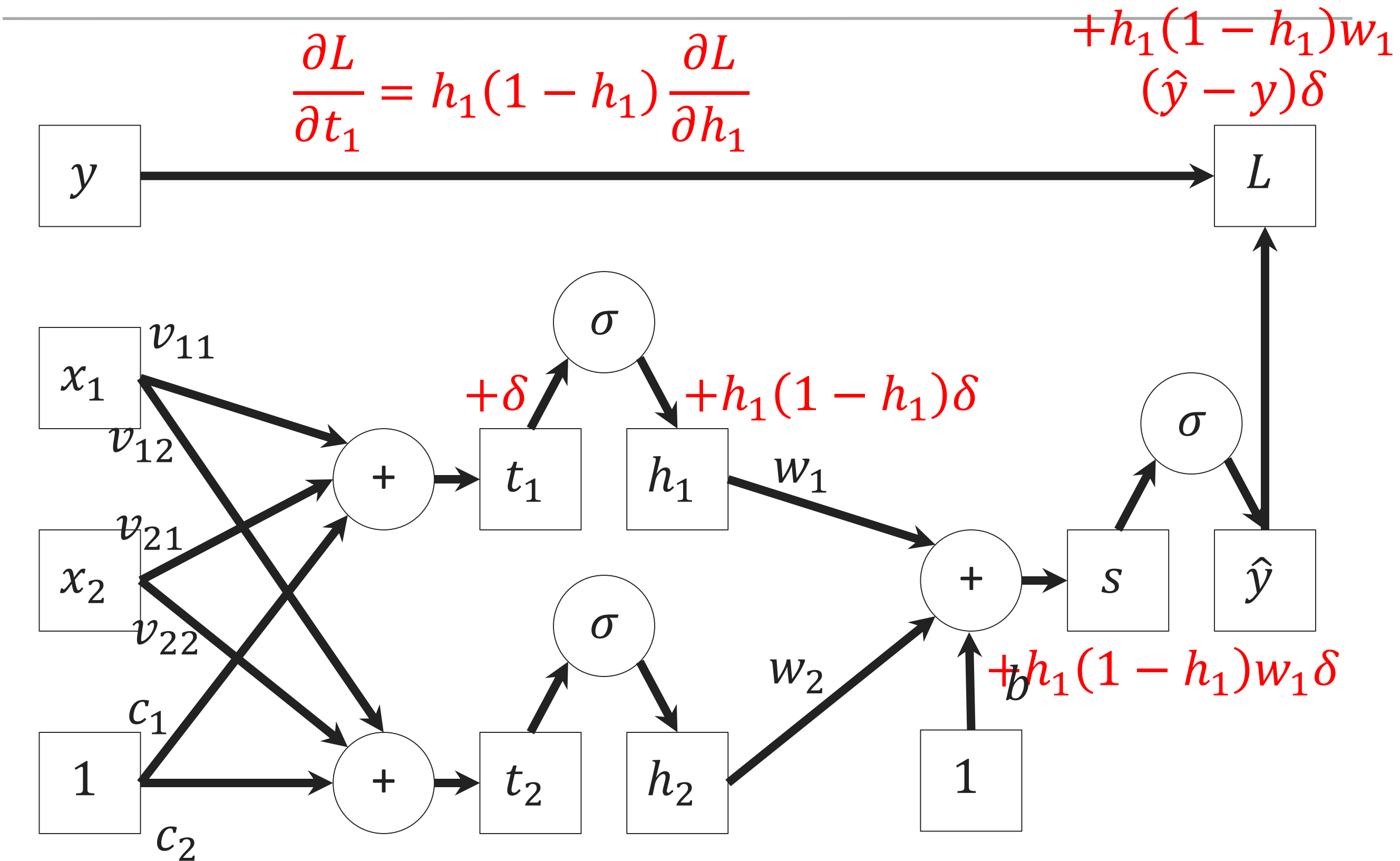
$$s = \mathbf{hw} + b = [h_1 \quad h_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b$$
$$\hat{y} = \sigma(s)$$

Forward pass

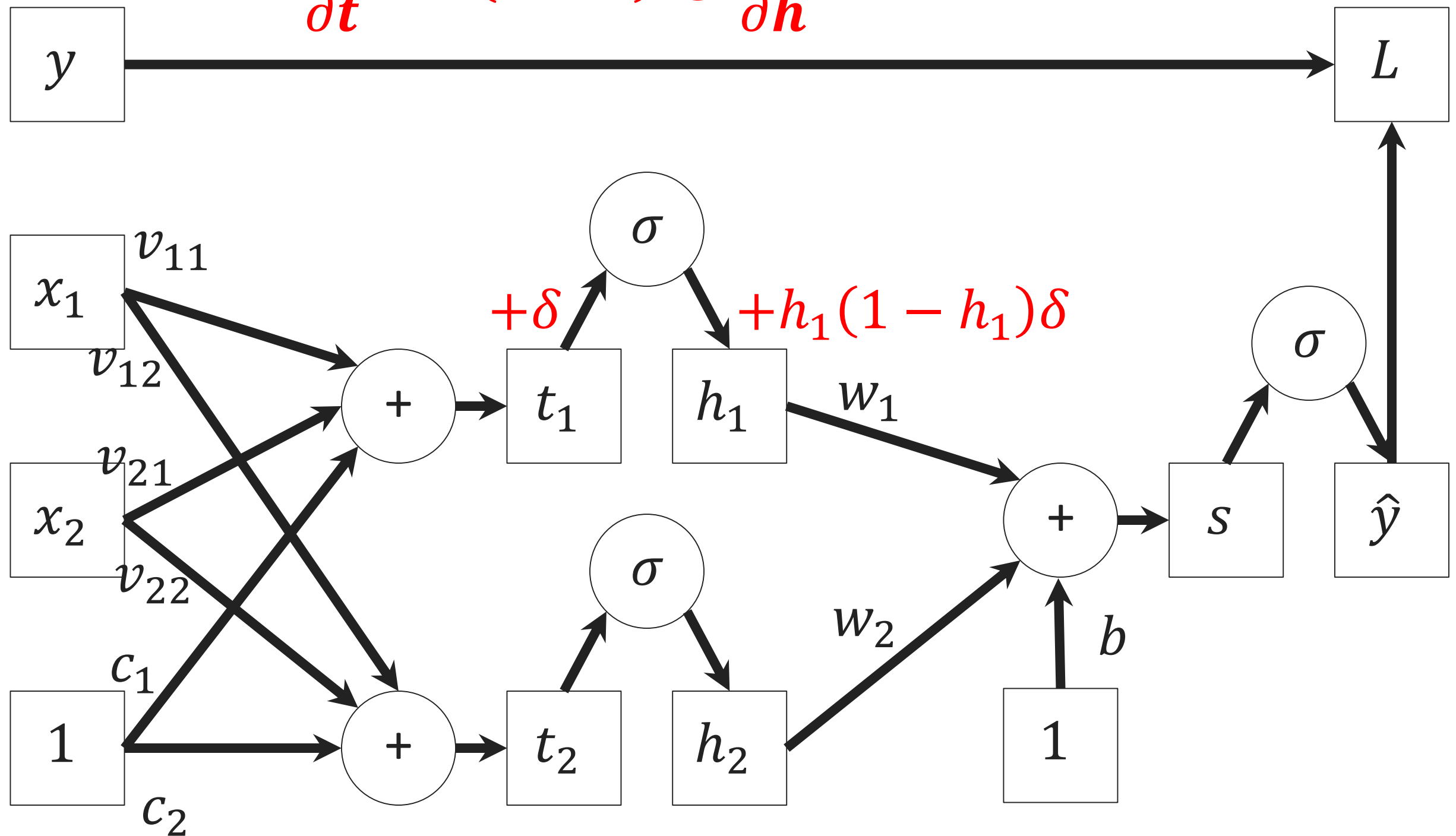


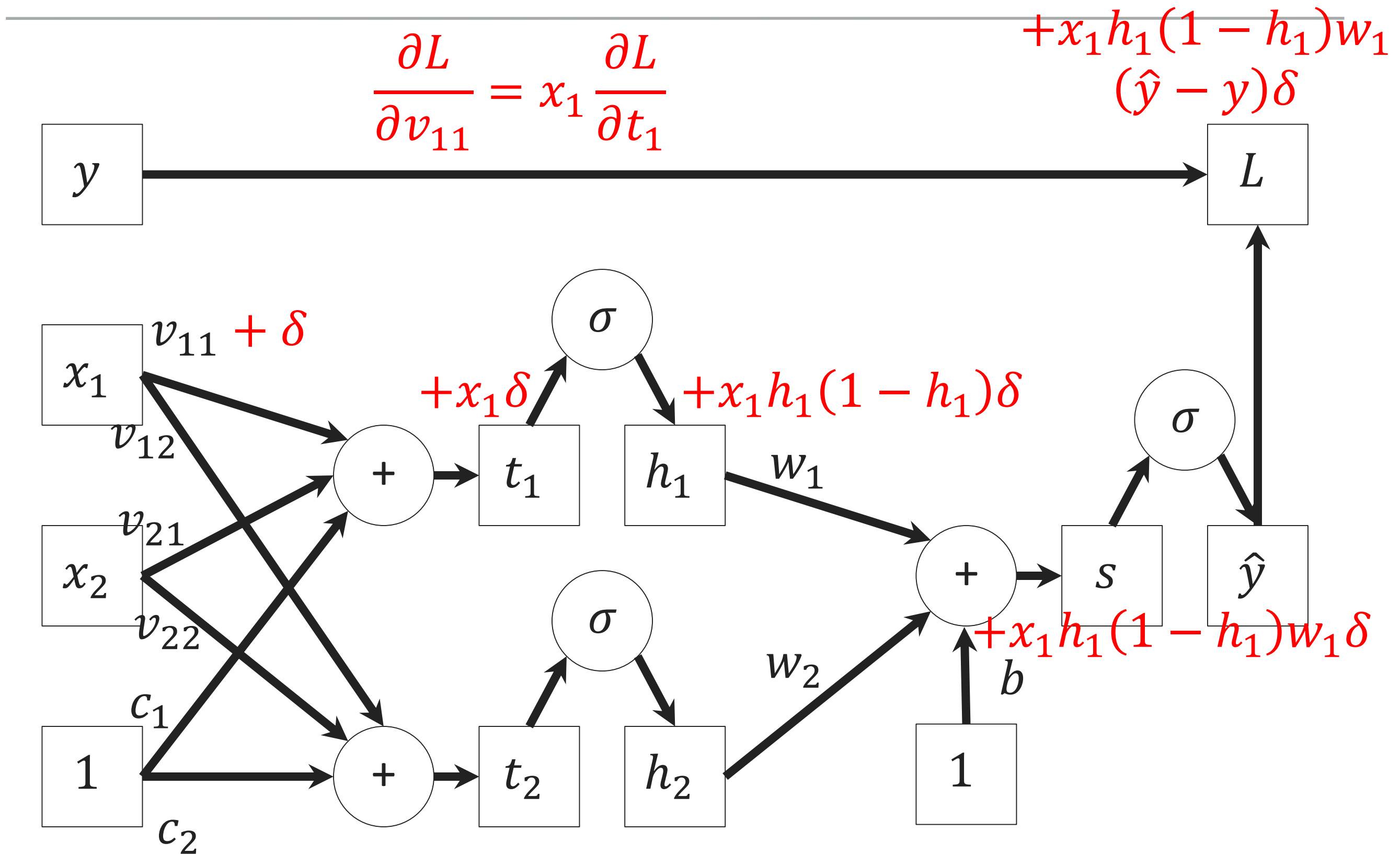




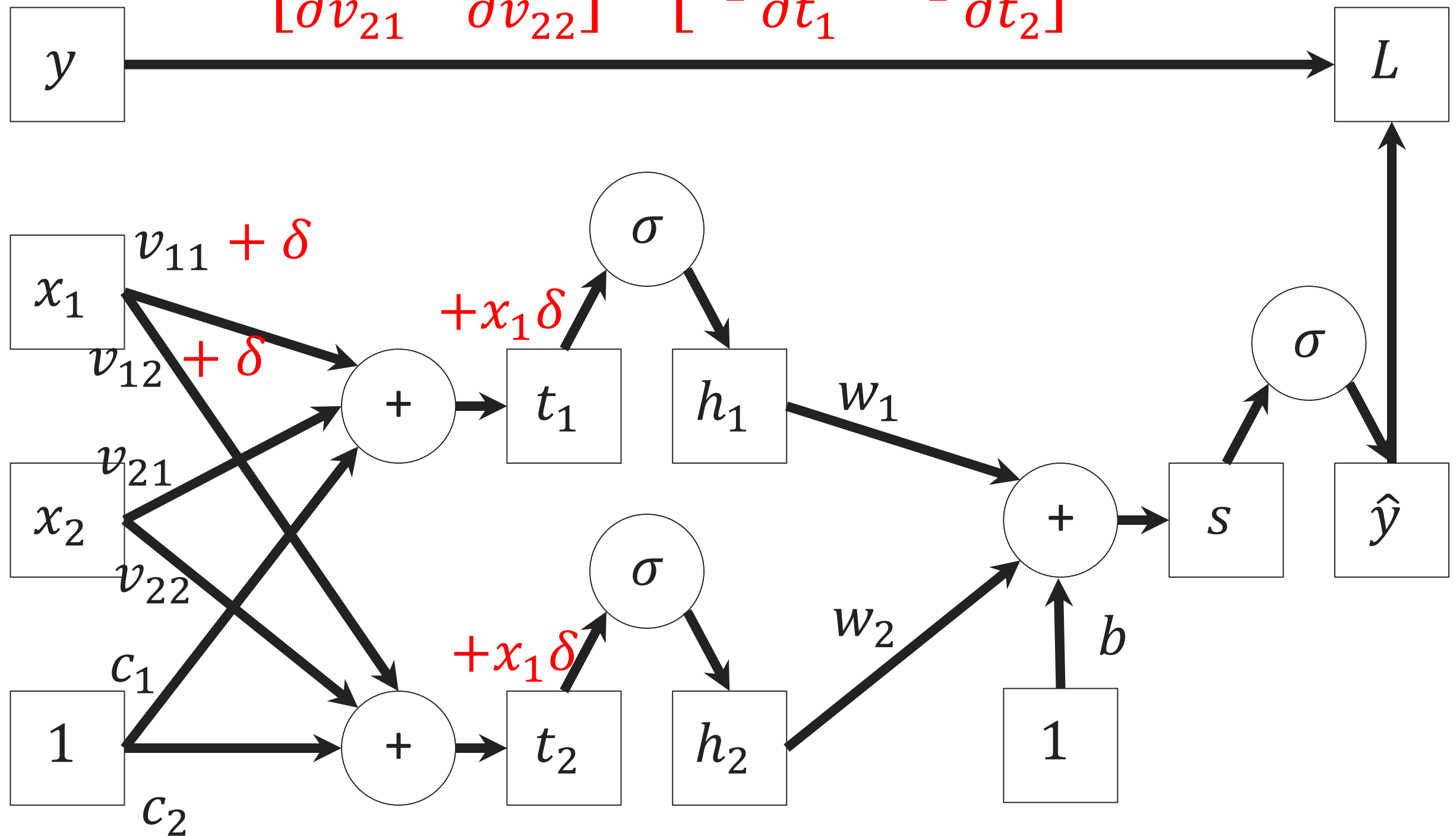


$$\frac{\partial L}{\partial t} = h(1 - h) \odot \frac{\partial L}{\partial h}$$

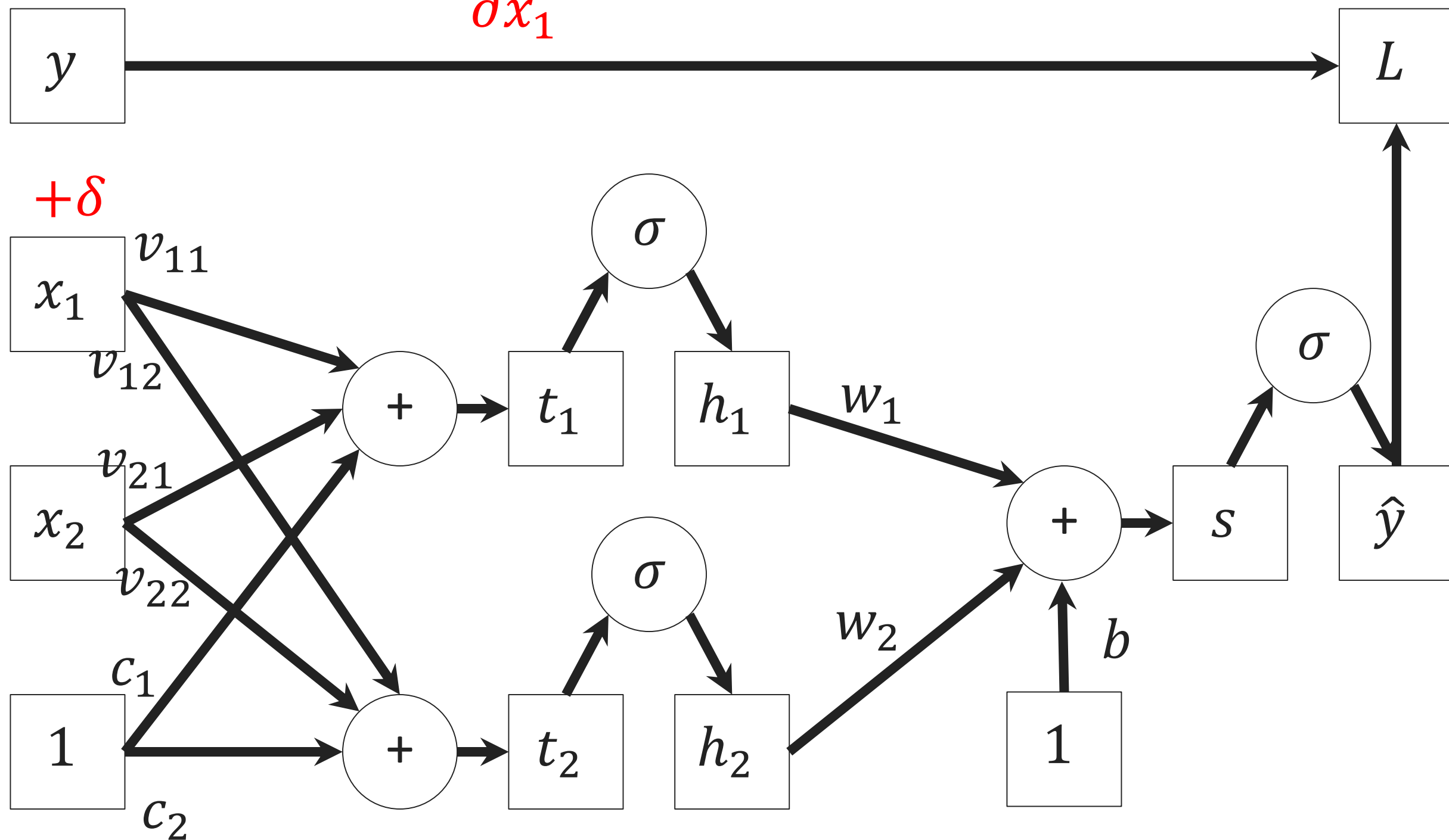




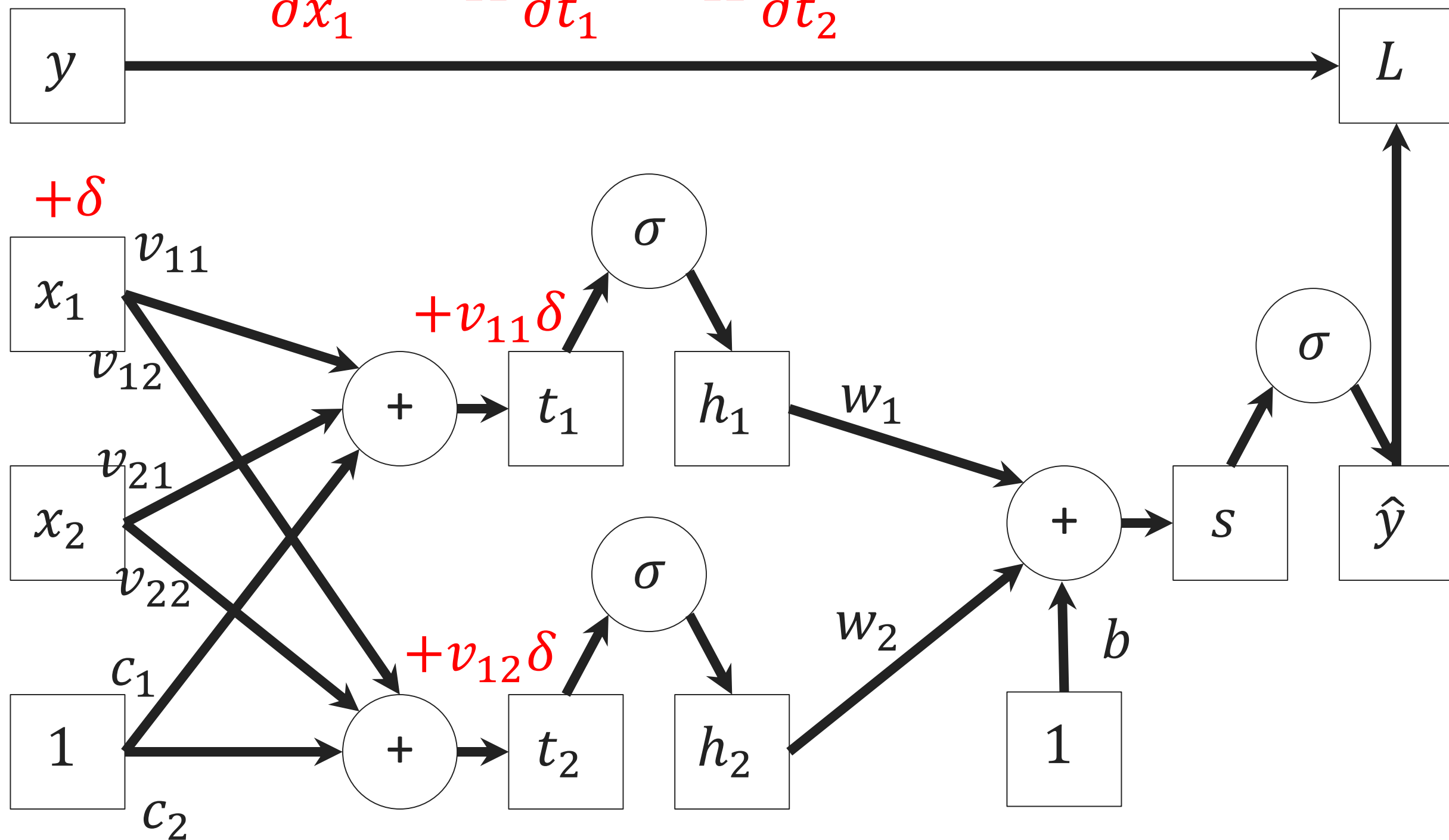
$$\frac{\partial L}{\partial \mathbf{V}} = \begin{bmatrix} \frac{\partial L}{\partial v_{11}} & \frac{\partial L}{\partial v_{12}} \\ \frac{\partial L}{\partial v_{21}} & \frac{\partial L}{\partial v_{22}} \end{bmatrix} = \begin{bmatrix} x_1 \frac{\partial L}{\partial t_1} & x_1 \frac{\partial L}{\partial t_2} \\ x_2 \frac{\partial L}{\partial t_1} & x_2 \frac{\partial L}{\partial t_2} \end{bmatrix} = \frac{\partial L}{\partial \mathbf{t}} \mathbf{x}^T$$



$$\frac{\partial L}{\partial x_1} = ?$$

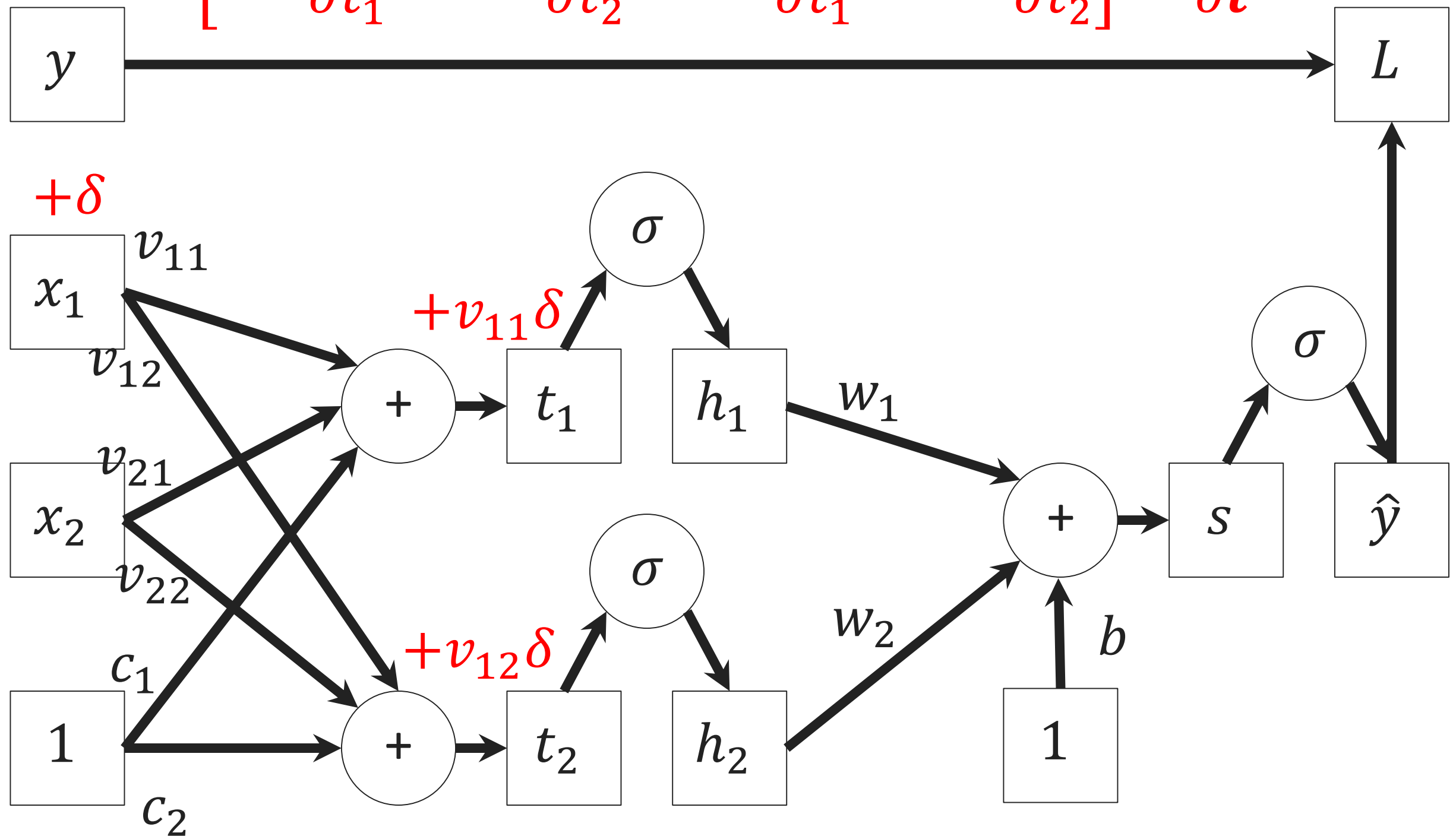


$$\frac{\partial L}{\partial x_1} = v_{11} \frac{\partial L}{\partial t_1} + v_{12} \frac{\partial L}{\partial t_2}$$

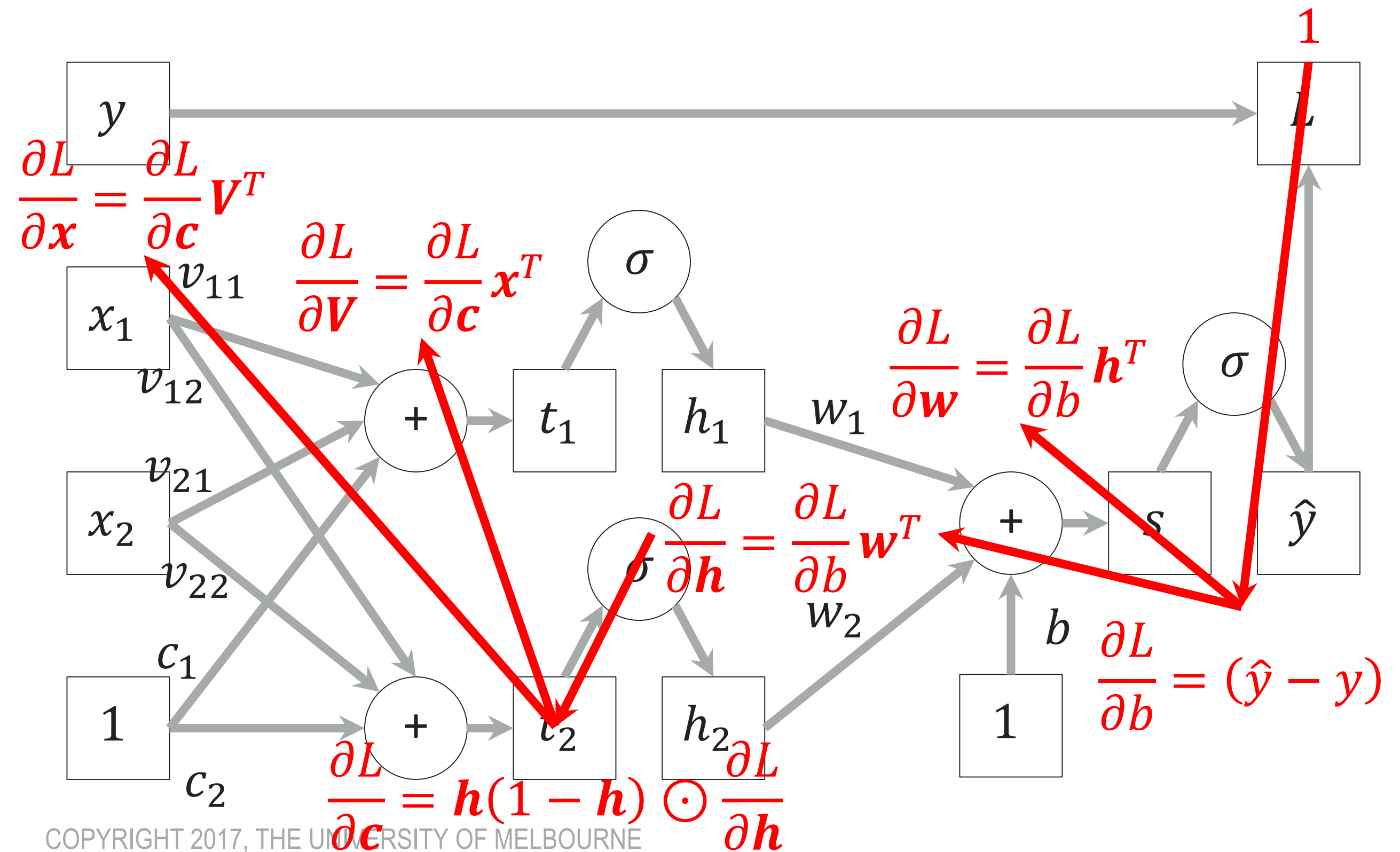


$$\frac{\partial L}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial L}{\partial x_1} & \frac{\partial L}{\partial x_2} \end{bmatrix}$$

$$= \begin{bmatrix} v_{11} \frac{\partial L}{\partial t_1} + v_{12} \frac{\partial L}{\partial t_2} & v_{21} \frac{\partial L}{\partial t_1} + v_{22} \frac{\partial L}{\partial t_2} \end{bmatrix} = \frac{\partial L}{\partial \mathbf{t}} \mathbf{v}^T$$



Backpropagation



Outline

- Review the lecture, background knowledge, etc.
 - Gradient descent & stochastic gradient descent (SGD)
 - Gradient and backpropagation
 - Logistic regression
 - Neural networks with one hidden layer

- Notebook tasks
 - Task 1: Multi-layer perceptron, SGD