

# The Mathematical Engineering of Deep Learning

Chapter 2 - Lecture 2 - Part (3/3)

---

B. Liquet<sup>1,2</sup> and S. Moka<sup>3</sup> and Y. Nazarathy<sup>3</sup>

<sup>1</sup> Macquarie University <sup>2</sup> LMAP, Université de Pau et des Pays de L'Adour <sup>3</sup> The University of Queensland

At the moment we use a lot of images, videos and data that are created by other authors. We have a duty to support the people who create it. We can do this by making it easier for them to share their work. We can do this by making it easier for them to share their work.

# Outline of Lecture

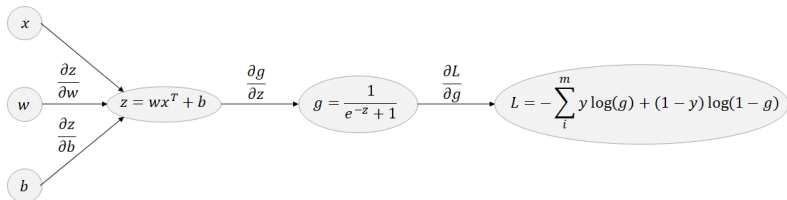
- Backpropagation
  - Chain rule
  - computational graph
- Implementation

# Backpropagation

## Backpropagation

The key tool adopted by the neural network community to update the weights. This method exploits the derivative with respect to each weight  $w$  using the **chain rule** (univariate and multivariate rules)

The chain rule is used and generally illustrated through a **computation graph**



## Chain rule

The univariate **chain rule** and the multivariate chain rule are the keys concept to calculate the derivative of cost with respect to any weight in the network.

- Univariate chain rule

$$\frac{\partial f(g(w))}{\partial w} = \frac{\partial f(g(w))}{\partial g(w)} \cdot \frac{\partial g(w)}{\partial w}$$

## Example: univariate chain rule

### Univariate logistic least squares model

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

### Computing the derivatives using chain rules

$$\frac{\partial \mathcal{L}}{\partial y} = y - t$$

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial y} \sigma'(z)$$

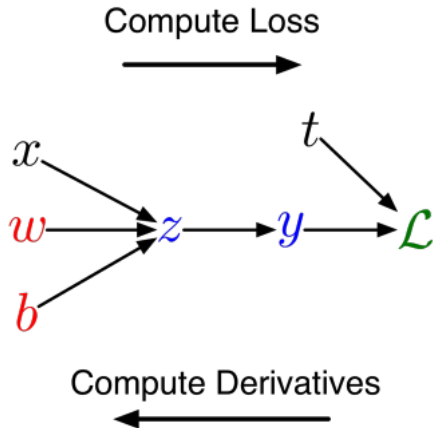
$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} x$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z}$$

**Goal:** write a program that efficiently computes the derivatives and not get closed form

## Computation graph.

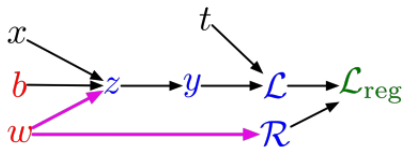
Help: Draw a diagram, **Computation graph**, to represent our computations



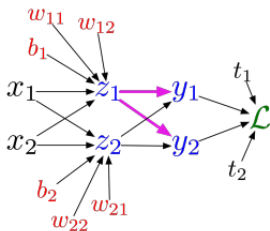
## New Challenge

Need **multivariate Chain Rule** when computation graph has fan-out  $> 1$ .

### $L_2$ -Regularized regression



### Multiclass logistic regression



# Multivariate Chain Rule

## Part I

Let  $z = f(x, y)$ ,  $x = g(t)$  and  $y = h(t)$ , where  $f, g$  and  $h$  are differentiable functions. Then  $z = f(x, y) = f(g(t), h(t))$  is a function of  $t$ , and

$$\begin{aligned}\frac{dz}{dt} = \frac{df}{dt} &= f_x(x, y) \frac{dx}{dt} + f_y(x, y) \frac{dy}{dt} \\ &= \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}.\end{aligned}$$



# Multivariate Chain Rule

## Part II

a) Let  $z = f(x, y)$ ,  $x = g(s, t)$  and  $y = h(s, t)$ , where  $f, g$  and  $h$  are differentiable functions. Then  $z$  is a function of  $s$  and  $t$ , and

$$\frac{\partial z}{\partial s} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial s}$$

and

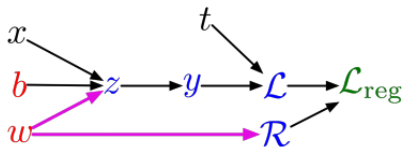
$$\frac{\partial z}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$$

b) Let  $z = f(x_1, x_2, \dots, x_m)$  be a differentiable function of  $m$  variables, where each of the  $x_i$  is a differentiable function of the variables  $t_1, t_2, \dots, t_n$ . Then  $z$  is a function of the  $t_i$ , and

$$\frac{\partial z}{\partial t_j} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t_j} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t_j} + \dots + \frac{\partial f}{\partial x_m} \frac{\partial x_m}{\partial t_j}.$$

## Case of Regularized regression

### $L_2$ -Regularized regression



Forward pass:

$$z = wx + b$$

$$y = \sigma(z)$$

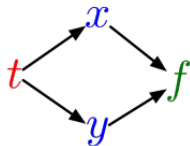
$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

$$R = \frac{1}{2}w^2$$

$$\mathcal{L}_{\text{reg}} = \mathcal{L} + \lambda R$$

## Multivariate Chain Rule: understand

- Consider a function  $f(x, y)$  and functions  $x(t)$  and  $y(t)$ . All the variables are scalar-valued.



$$\frac{d}{dt}f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

# Multivariate Chain Rule: understand

- Example

$$f(x, y) = y + e^{xy}$$

$$x(t) = \cos t$$

$$y(t) = t^2$$

- Using the chain rule

$$\begin{aligned}\frac{df}{dt} &= \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} \\ &= (ye^{xy}) \cdot (-\sin t) + (1 + xe^{xy}) \cdot 2t\end{aligned}$$

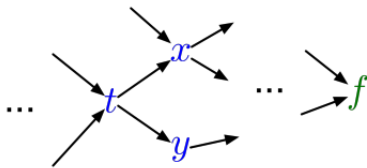
# Multivariate Chain Rule: general example

- Backpropagation framework

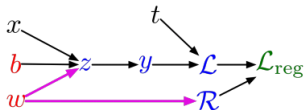
Mathematical expressions  
to be evaluated

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

Values already computed  
by our program



## $L_2$ -Regularized regression



### Forward pass

$$z = wx + b \rightarrow y = \sigma(z) \rightarrow \mathcal{L} = \frac{1}{2}(y - t)^2 \rightarrow \mathcal{R} = \frac{1}{2}w^2 \rightarrow \mathcal{L}_{reg} = \mathcal{L} + \lambda\mathcal{R}$$

### Backward pass

$$\frac{\partial \mathcal{L}_{reg}}{\partial \mathcal{R}} = \lambda \rightarrow \frac{\partial \mathcal{L}_{reg}}{\partial \mathcal{L}} = 1 \rightarrow \frac{\partial \mathcal{L}_{reg}}{\partial y} = \frac{\partial \mathcal{L}_{reg}}{\partial \mathcal{L}} \times \frac{\partial \mathcal{L}}{\partial y} = 1 \times (y - t)$$

$$\rightarrow \frac{\partial \mathcal{L}_{reg}}{\partial z} = \frac{\partial \mathcal{L}_{reg}}{\partial y} \times \frac{\partial y}{\partial z} = \frac{\partial \mathcal{L}_{reg}}{\partial y} \times \sigma'(z)$$

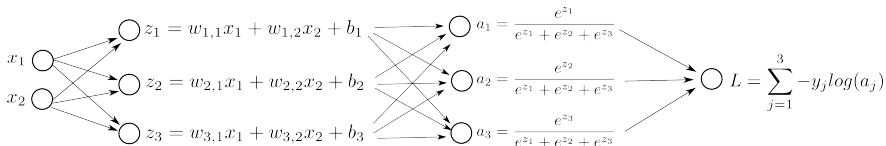
$$\rightarrow \frac{\partial \mathcal{L}_{reg}}{\partial w} = \frac{\partial \mathcal{L}_{reg}}{\partial z} \times \frac{\partial z}{\partial w} + \frac{\partial \mathcal{L}_{reg}}{\partial \mathcal{R}} \times \frac{\partial \mathcal{R}}{\partial w} = \frac{\partial \mathcal{L}_{reg}}{\partial z} x + \frac{\partial \mathcal{L}_{reg}}{\partial \mathcal{R}} w$$

$$\rightarrow \frac{\partial \mathcal{L}_{reg}}{\partial b} = \frac{\partial \mathcal{L}_{reg}}{\partial z} \times \frac{\partial z}{\partial b} = \frac{\partial \mathcal{L}_{reg}}{\partial z}$$

## Example: softmax neural network

Let consider a simple example with

- $K = 3$  ( $y \in \{1, 2, 3\}$ ) and two features ( $x_1$  and  $x_2$ ).
- $w_j \in \mathbb{R}^2$  and  $b_j \in \mathbb{R}$ ,  $j = 1, 2, 3$ .



Let's write as an example for  $\frac{\partial L}{\partial w_{2,1}}$ :

$$\begin{aligned} \frac{\partial L}{\partial w_{2,1}} &= \sum_{i=1}^3 \left( \frac{\partial L}{\partial a_i} \right) \left( \frac{\partial a_i}{\partial z_2} \right) \left( \frac{\partial z_2}{\partial w_{2,1}} \right) \\ &= \left( \frac{\partial L}{\partial a_1} \right) \left( \frac{\partial a_1}{\partial z_2} \right) \left( \frac{\partial z_2}{\partial w_{2,1}} \right) + \left( \frac{\partial L}{\partial a_2} \right) \left( \frac{\partial a_2}{\partial z_2} \right) \left( \frac{\partial z_2}{\partial w_{2,1}} \right) + \left( \frac{\partial L}{\partial a_3} \right) \left( \frac{\partial a_3}{\partial z_2} \right) \left( \frac{\partial z_2}{\partial w_{2,1}} \right) \end{aligned}$$

When we change  $w_{2,1}$ ;  $a_1$ ,  $a_2$  and  $a_3$  changes as a result. Then, the change of  $a_1$ ,  $a_2$  and  $a_3$  affects  $L$ . We sum up all the changes  $w_{2,1}$  produced over  $a_1$ ,  $a_2$  and

# Updating weights using Backpropagation

For neural network framework, the weight is updated using gradient descent concepts

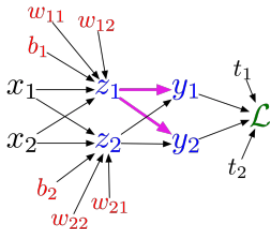
$$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

The main steps for updating weights are

1. Take a batch of training sample
2. Forward propagation to get the corresponding cost
3. Backpropagate the cost to get the gradients
4. update the weights using the gradients
5. Repeat step 1 to 4 for a number of iterations



## Multiclass logistic regression



$$z_\ell = \sum_j w_{\ell j} x_j + b_\ell$$

$$y_k = \frac{e^{z_k}}{\sum_\ell e^{z_\ell}}$$

$$\mathcal{L} = - \sum_k t_k \times \ln y_k$$

**TO DO:** Forward pass and Backward pass

# Take Home Message

- **Forward propagation:** input  $x$  is propagated through the network (graph) to reach the output  $\hat{y}$ . Then the cost could be computed  $\mathcal{J}(\theta)$
- **Back-propagation:** process to compute the gradient by using repeatedly the chain rule coupled with a memory that allows to reuse computations.