

VII.3] Backpropagation and The Chain Rule

When training a neural network, most of the time is being spent on computing the gradient of the loss function with respect to the many parameters of the neural net.

Finding an explicit formula for the gradient and writing the code to evaluate the gradient is not reasonable for such a large optimization problem.

There are two alternatives to finding the gradient explicitly.

One alternative is to use finite differences to approximate all the partial derivatives in the gradient:

$$\nabla f(x)_i = \frac{\partial f}{\partial x_i} \approx \frac{f(x + h e_i) - f(x)}{h},$$

where e_i is the i^{th} column of the identity matrix. However, this would require us to evaluate the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ many times, and we would only get an approximation.

A better approach is to use AD (automatic differentiation) which produces a routine that will exactly evaluate the derivatives.

Two ways AD is implemented are:

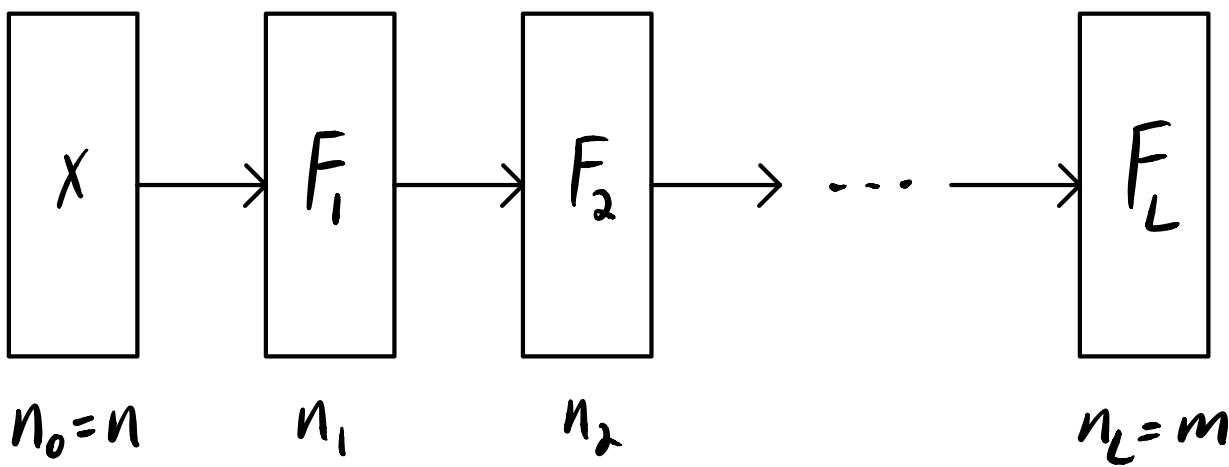
- (1) source code transformation (SCT),
- (2) operator overloading (OO).

The Julia package Zygote.jl is an SCT implementation while the ForwardDiff.jl package uses OO.

Forward-mode AD

Suppose $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is given by a composition of intermediate calculations:

$$\begin{aligned} F(x) &= (F_L \circ \dots \circ F_2 \circ F_1)(x) \\ &= F_L(\dots F_2(F_1(x))). \end{aligned}$$



Here $F_i: \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}$, for $i=1, \dots, L$.

We use $F'(x) = \frac{\partial F}{\partial x}(x)$ to represent the $m \times n$ Jacobian matrix of $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$:

$$F(x+d) \approx F(x) + F'(x)d, \quad \forall d \in \mathbb{R}^n \text{ small.}$$

Example: $F: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ $F'(x)$ is 3×2

$$F(x) = \begin{bmatrix} x_1 + x_2 \\ x_1 x_2 \\ x_1^2 - x_2 \end{bmatrix} \quad F'(x) = \begin{bmatrix} 1 & 1 \\ x_2 & x_1 \\ 2x_1 & -1 \end{bmatrix}$$

//

We can use the chain rule to take the derivative of a composition:

$$F(x) = F_2(F_1(x)) \Rightarrow F'(x) = F'_2(F_1(x)) \cdot F'_1(x).$$

With $F_1: \mathbb{R}^n \rightarrow \mathbb{R}^{n_1}$, $F_2: \mathbb{R}^{n_1} \rightarrow \mathbb{R}^m$, and $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, we have:

$$F'(x) = F'_2(F_1(x)) \cdot F'_1(x).$$

$m \times n$

$m \times n_1$

$n_1 \times n$

Example: $F_1: \mathbb{R}^2 \rightarrow \mathbb{R}^2$, $F_2: \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$F_1(x) = \begin{bmatrix} x_1, -x_2 \\ x_1 x_2 \end{bmatrix}, \quad F_2(x) = \begin{bmatrix} x_1^2 \\ x_1 - x_2 \end{bmatrix}$$

$$\Rightarrow F(x) = F_2(F_1(x)) = \begin{bmatrix} (x_1 - x_2)^2 \\ x_1 - x_2 - x_1 x_2 \end{bmatrix}.$$

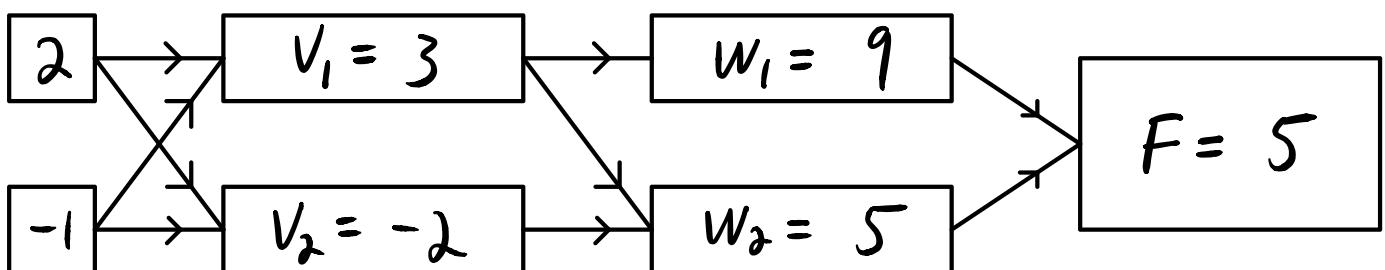
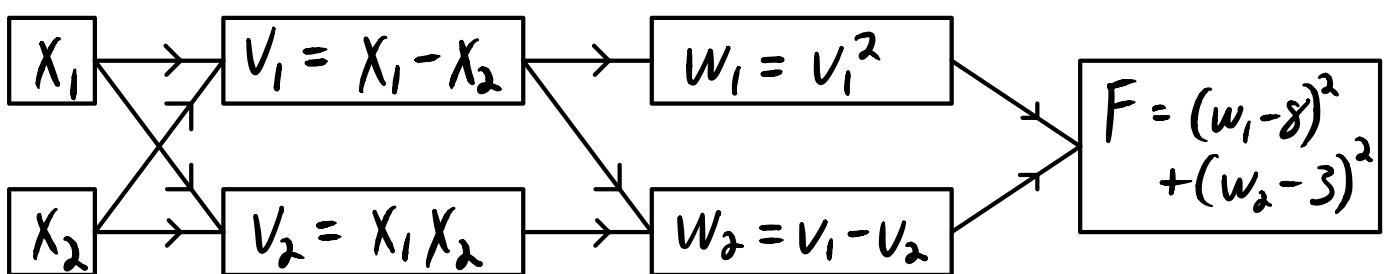
$$F'_1(x) = \begin{bmatrix} 1 & -1 \\ x_2 & x_1 \end{bmatrix}, \quad F'_2(x) = \begin{bmatrix} 2x_1 & 0 \\ 1 & -1 \end{bmatrix}.$$

$$F'(x) = \begin{bmatrix} 2(x_1 - x_2) & -2(x_1 - x_2) \\ 1 - x_2 & -1 - x_1 \end{bmatrix}$$

$$\begin{aligned} F_2'(F_1(x)) \cdot F_1'(x) &= \begin{bmatrix} 2(x_1 - x_2) & 0 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ x_2 & x_1 \end{bmatrix} \\ &= \begin{bmatrix} 2(x_1 - x_2) & -2(x_1 - x_2) \\ 1 - x_2 & -1 - x_1 \end{bmatrix}. \quad \checkmark \end{aligned}$$

Suppose $F_3(x) = (x_1 - 8)^2 + (x_2 - 3)^2$ and

$$F(x) = F_3(F_2(F_1(x))).$$



To compute F , we pass forward through the network. Now we want to compute

$$F(x+d) \approx F(x) + \nabla F(x)^T d$$

$$\begin{aligned} F'(x) &= \begin{bmatrix} \frac{\partial F}{\partial x_1} & \frac{\partial F}{\partial x_2} \end{bmatrix} \quad (2 \times 1) \\ F: \mathbb{R}^2 &\rightarrow \mathbb{R}^1 \\ &= \nabla F(x)^T. \end{aligned}$$

First we compute $\frac{\partial F}{\partial x_1}$ with a forward

pass:

$$\begin{aligned} x_1 &= 2 \\ x_2 &= -1 \end{aligned}$$

$$v_1 = x_1 - x_2$$

$$v_2 = x_1 x_2$$

$$\frac{\partial x_1}{\partial x_1} = 1 \quad \frac{\partial v_1}{\partial x_1} = \frac{\partial x_1}{\partial x_1} - \frac{\partial x_2}{\partial x_1} = 1$$

...

$$\frac{\partial x_2}{\partial x_1} = 0 \quad \frac{\partial v_2}{\partial x_1} = \frac{\partial x_1}{\partial x_1} x_2 + x_1 \frac{\partial x_2}{\partial x_1} = -1$$

$$\begin{aligned} 1 \cdot (-1) + 2 \cdot 0 \\ = -1 \end{aligned}$$

$$v_1 = 3$$

$$v_2 = -2$$

$$2 \cdot 3 \cdot 1 = 6$$

$$w_1 = v_1^2$$

$$w_2 = v_1 - v_2$$

...

$$\frac{\partial v_1}{\partial x_1} = 1 \quad \frac{\partial w_1}{\partial x_1} = 2 v_1 \frac{\partial v_1}{\partial x_1} = 6$$

...

$$\frac{\partial v_2}{\partial x_1} = -1 \quad \frac{\partial w_2}{\partial x_1} = \frac{\partial v_1}{\partial x_1} - \frac{\partial v_2}{\partial x_1} = 2$$

$$1 - (-1) = 2$$

$$w_1 = 9$$

$$w_2 = 5$$

...

$$\frac{\partial w_1}{\partial x_1} = 6$$

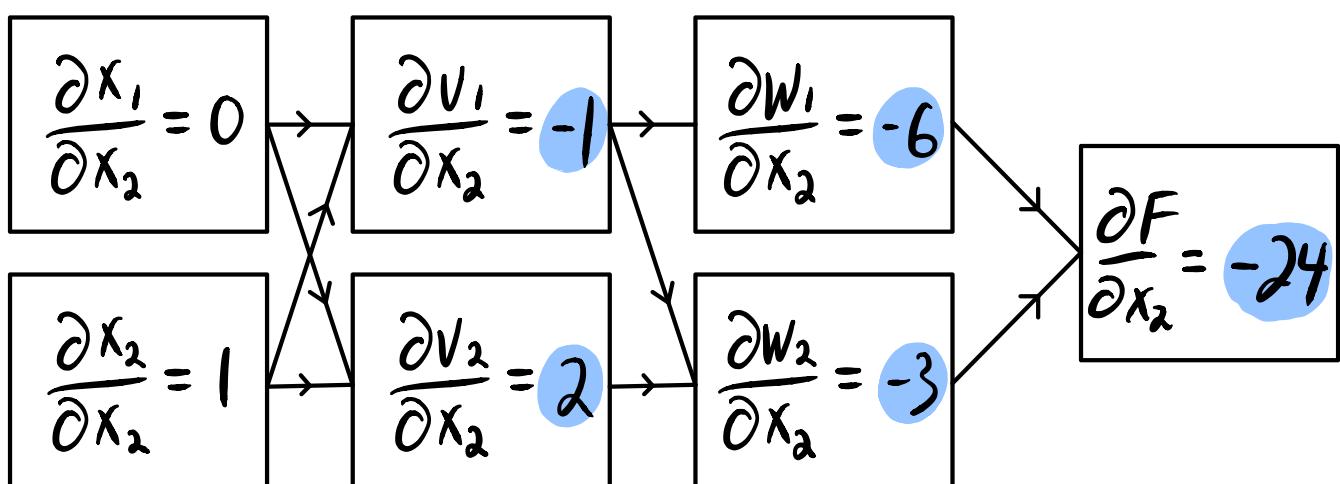
$$\frac{\partial w_2}{\partial x_1} = 2$$

$$F = (w_1 - 8)^2 + (w_2 - 3)^2$$

$$\begin{aligned} \frac{\partial F}{\partial x_1} &= 2(w_1 - 8) \frac{\partial w_1}{\partial x_1} \\ &\quad + 2(w_2 - 3) \frac{\partial w_2}{\partial x_1} = 20 \end{aligned}$$

$$2 \cdot (9-8) \cdot 6 + 2 \cdot (5-3) \cdot 2 = 12 + 8 = 20$$

Computing $\frac{\partial F}{\partial x_2}$ requires another forward pass through the network:



Thus, forward-mode AD requires a separate pass through the network for each variable. This will be a lot of work if there are many variables, like in neural nets.

Let

$$v = F_1(x) = \begin{bmatrix} x_1 - x_2 \\ x_1 x_2 \end{bmatrix}, \quad w = F_2(v) = \begin{bmatrix} v_1^2 \\ v_1 - v_2 \end{bmatrix},$$

and $F = F_3(w) = (w_1 - 8)^2 + (w_2 - 3)^2$.

Then

$$F'(x) = F_3'(w) \cdot F_2'(v) \cdot F_1'(x)$$

$$= [2(w_1 - 8) \quad 2(w_2 - 3)] \begin{bmatrix} 2v_1 & 0 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ x_2 & x_1 \end{bmatrix}$$

$$= [2 \quad 4] \left(\begin{bmatrix} 6 & 0 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \right)$$

$$= [2 \quad 4] \begin{bmatrix} 6 & -6 \\ 2 & -3 \end{bmatrix} \quad \text{first pass}$$

second pass

$$= [20 \quad -24].$$

Forward-mode AD \equiv mult. from the right.

We could compute $F'(x)$ more efficiently by multiplying from the left:

$$\begin{aligned} F'(x) &= \left(\begin{bmatrix} 2 & 4 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 6 & 0 \\ -1 & 2 \end{bmatrix} \right) \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \\ &= \begin{bmatrix} 16 & -4 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \\ &= \begin{bmatrix} 20 & -24 \end{bmatrix} \end{aligned}$$

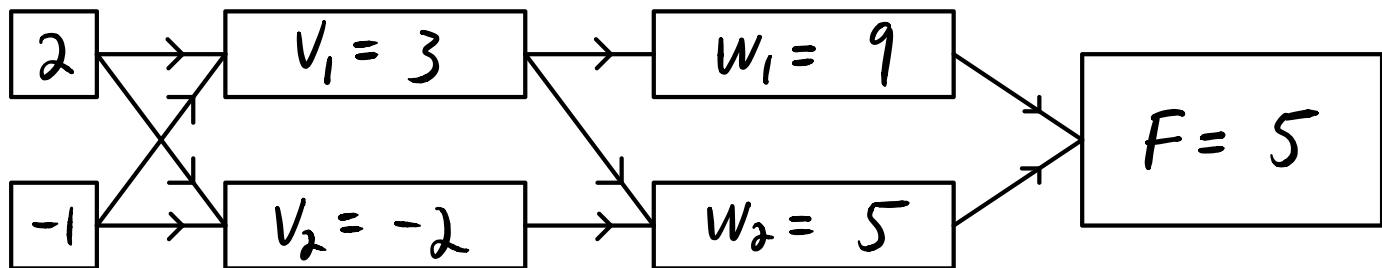
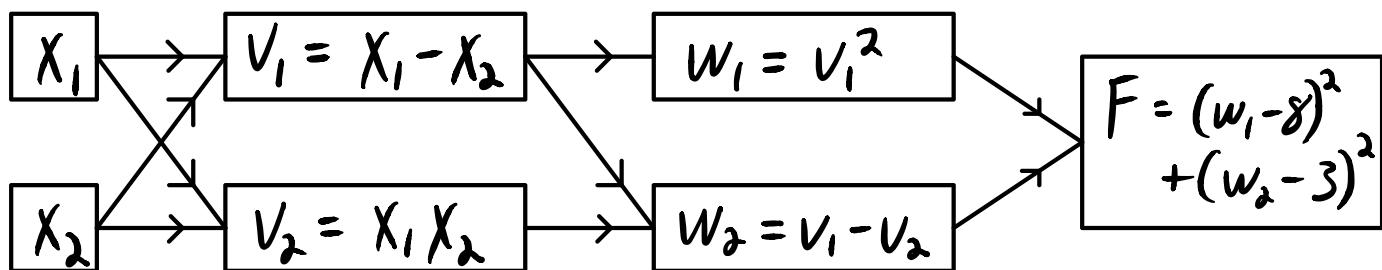
Notice that we only needed to do vector-matrix multiplication; we avoided doing the matrix-matrix multiplication. This can save us from doing a lot of extra work.

Mult. from the left \equiv Backward-mode AD

Backward-mode AD

This mode is also called reverse-mode AD or backpropagation.

We still begin with a forward pass to compute $F(x)$.



Then we compute all partial derivatives with a single backward pass.

$$\frac{\partial F}{\partial F} = 1$$

$\frac{\partial F}{\partial w_1} = 2(w_1 - 8) = 2$

$\frac{\partial F}{\partial v_1} \dots$

$\frac{\partial F}{\partial w_2} = 2(w_2 - 3) = 4$

$\frac{\partial F}{\partial v_2} \dots$

$$\dots = \frac{\partial F}{\partial w_1} \cdot \frac{\partial w_1}{\partial v_1} + \frac{\partial F}{\partial w_2} \cdot \frac{\partial w_2}{\partial v_1} = 16$$

$\frac{\partial F}{\partial x_1} \dots$

$$\dots = \frac{\partial F}{\partial w_2} \cdot \frac{\partial w_2}{\partial v_2} = -4$$

$\frac{\partial F}{\partial x_2} \dots$

$$\dots = \frac{\partial F}{\partial v_1} \cdot \frac{\partial v_1}{\partial x_1} + \frac{\partial F}{\partial v_2} \cdot \frac{\partial v_2}{\partial x_1} = 20$$

Single pass

$$\dots = \frac{\partial F}{\partial v_1} \cdot \frac{\partial v_1}{\partial x_2} + \frac{\partial F}{\partial v_2} \cdot \frac{\partial v_2}{\partial x_2} = -24$$

$$F'(x) = \left(\begin{bmatrix} 2 & 4 \end{bmatrix} \begin{bmatrix} 6 & 0 \\ 1 & -1 \end{bmatrix} \right) \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 16 & -4 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 20 & -24 \end{bmatrix}$$

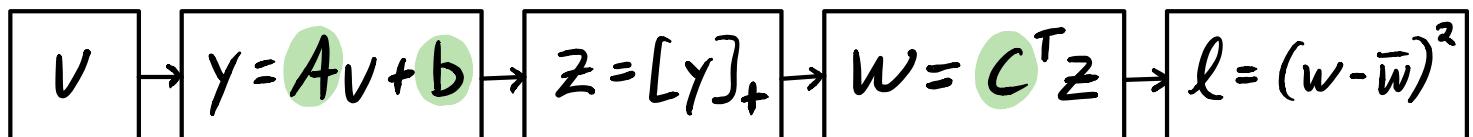
Note: If $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, then using forward-mode AD to compute $F'(x)$ requires n forward passes through the computational graph, while backward-mode AD requires m backward passes. Therefore, if $n < m$, forward-mode should be used and if $n > m$, backward-mode is preferred. Training a neural net requires minimizing a loss function $F: \mathbb{R}^n \rightarrow \mathbb{R}$ where n , the number of parameters, is very large; thus, backward-mode AD will be far more efficient than forward-mode AD.

Backpropagation on a neural net

Example:

$$A = \begin{bmatrix} 1 & 0 \\ -2 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad c = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$X = (A, b, c)$ are the parameters



Data: $V = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad \bar{W} = 1$

$$V = \begin{bmatrix} 2 \\ -1 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 3 \\ -5 \end{bmatrix} \rightarrow Z = \begin{bmatrix} 3 \\ 0 \end{bmatrix} \rightarrow W = 3 \rightarrow l = 4$$

$$\frac{\partial l}{\partial w} = 2(w - \bar{w}) = 4$$

$$(l = (W - \bar{W})^2)$$

$$\frac{\partial l}{\partial c_1} = \frac{\partial l}{\partial w} \cdot \frac{\partial w}{\partial c_1} = 4z_1 = 12$$

$$(W = c_1 z_1 + c_2 z_2)$$

$$\frac{\partial l}{\partial c_2} = \frac{\partial l}{\partial w} \cdot \frac{\partial w}{\partial c_2} = 4z_2 = 0$$

$$\frac{\partial l}{\partial z_1} = \frac{\partial l}{\partial w} \cdot \frac{\partial w}{\partial z_1} = 4c_1 = 4$$

$$(W = c_1 z_1 + c_2 z_2)$$

$$\frac{\partial l}{\partial z_2} = \frac{\partial l}{\partial w} \cdot \frac{\partial w}{\partial z_2} = 4c_2 = 4$$

$$\frac{\partial l}{\partial y_1} = \frac{\partial l}{\partial z_1} \cdot \frac{\partial z_1}{\partial y_1} = 4 \cdot 1 = 4$$

$$(z_1 = y_1)$$

$$\frac{\partial l}{\partial y_2} = \frac{\partial l}{\partial z_2} \cdot \frac{\partial z_2}{\partial y_2} = 4 \cdot 0 = 0$$

$$(z_2 = 0)$$

$$\frac{\partial l}{\partial a_{11}} = \frac{\partial l}{\partial y_1} \cdot \frac{\partial y_1}{\partial a_{11}} = 4v_1 = 8$$

$$\frac{\partial l}{\partial a_{12}} = \frac{\partial l}{\partial y_1} \cdot \frac{\partial y_1}{\partial a_{12}} = 4v_2 = -4$$

$$\frac{\partial l}{\partial b_1} = \frac{\partial l}{\partial y_1} \cdot \frac{\partial y_1}{\partial b_1} = 4 \cdot 1 = 4$$

$$(y_1 = a_{11}v_1 + a_{12}v_2 + b_1)$$

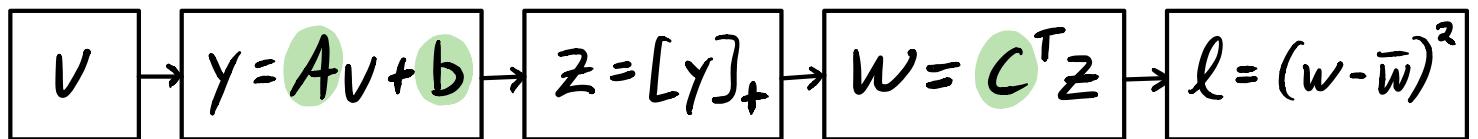
$$\frac{\partial l}{\partial a_{21}} = \frac{\partial l}{\partial y_2} \cdot \frac{\partial y_2}{\partial a_{21}} = 0 \cdot v_1 = 0$$

$$(y_2 = a_{21}v_1 + a_{22}v_2 + b_2)$$

$$\frac{\partial l}{\partial b_2} = \frac{\partial l}{\partial y_2} \cdot \frac{\partial y_2}{\partial b} = 0 \cdot 1 = 0$$

$$\therefore \nabla l(A, b, c) = \left(\begin{bmatrix} 8 & -4 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 4 \\ 0 \end{bmatrix}, \begin{bmatrix} 12 \\ 0 \end{bmatrix} \right)$$

Using matrices to compute ∇l :



$$\frac{\partial \ell}{\partial c}(\Delta c) = \frac{\partial \ell}{\partial w} \cdot \frac{\partial w}{\partial c}(c)$$

$$w + \Delta w = (c + \Delta c)^T z$$

$$\Delta w = z^T \Delta c$$

$$= 2(w - \bar{w}) z^T \Delta c = 4 [3 \ 0] \Delta c$$

$$= [12 \ 0] \Delta c = \left\langle \begin{bmatrix} 12 \\ 0 \end{bmatrix}, \Delta c \right\rangle \therefore \nabla_c \ell = \begin{bmatrix} 12 \\ 0 \end{bmatrix}$$

$$\frac{\partial \ell}{\partial A}(\Delta A) = \frac{\partial \ell}{\partial w} \cdot \frac{\partial w}{\partial Z} \cdot \frac{\partial Z}{\partial Y} \cdot \frac{\partial Y}{\partial A}(\Delta A)$$

$$y + \Delta y = (A + \Delta A)V + b$$

$$\Delta y = \Delta A V$$

$$= 2(w - \bar{w}) [c_1 \ c_2] \begin{bmatrix} [y_1]'_+ \\ [y_2]'_+ \end{bmatrix} (\Delta A V)$$

$$= 4 [1 \ 1] \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} (\Delta A V) = [4 \ 4] \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} (\Delta A V)$$

$$= [4 \ 0] \Delta A \begin{bmatrix} 2 \\ -1 \end{bmatrix} = [4 \Delta a_{11} \ 4 \Delta a_{12}] \begin{bmatrix} 2 \\ -1 \end{bmatrix}$$

$$= 8 \Delta a_{11} - 4 \Delta a_{12} = \left\langle \begin{bmatrix} 8 & -4 \\ 0 & 0 \end{bmatrix}, \Delta A \right\rangle$$

$$\therefore \nabla_A \ell = \begin{bmatrix} 8 & -4 \\ 0 & 0 \end{bmatrix}.$$

$$\frac{\partial \ell}{\partial b}(\Delta b) = \frac{\partial \ell}{\partial w} \cdot \frac{\partial w}{\partial z} \cdot \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial b}(\Delta b)$$

$$Y + \Delta y = Av + b + \Delta b$$

$$\Delta y = \Delta b$$

$$= 4 [1 \ 1] \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \Delta b = [4 \ 4] \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \Delta b$$

$$= [4 \ 0] \Delta b = \left\langle \begin{bmatrix} 4 \\ 0 \end{bmatrix}, \Delta b \right\rangle$$

$$\therefore \nabla_b \ell = \begin{bmatrix} 4 \\ 0 \end{bmatrix}.$$

$$\therefore \nabla \ell(A, b, c) = \left(\begin{bmatrix} 8 & -4 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 4 \\ 0 \end{bmatrix}, \begin{bmatrix} 12 \\ 0 \end{bmatrix} \right)$$

$$l: \mathbb{R} \rightarrow \mathbb{R} \quad f: \mathbb{R}^{2 \times 2} \times \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$w: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$z: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$y: \mathbb{R}^{2 \times 2} \times \mathbb{R}^2 \rightarrow \mathbb{R}^2 \quad \frac{\partial y}{\partial A}: \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}^2$$

$$\frac{\partial y}{\partial A}(A) = \begin{bmatrix} \langle \begin{bmatrix} v_1 & v_2 \\ 0 & 0 \end{bmatrix}, A \rangle \\ \langle \begin{bmatrix} 0 & 0 \\ v_1 & v_2 \end{bmatrix}, A \rangle \end{bmatrix} = Av$$