**1.** $\quad w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)}) \cdots \cdots$ (1)

$$\underset{w}{\arg\min} \; f(w^{(t)}) + \langle w - w^{(t)}, \nabla f(w^{(t)}) \rangle + \frac{\lambda}{2} \| w - w^{(t)} \|^2 \quad (4)$$

Take derivative of (4) w.r.t. $w$, & set to $0$.

$$0 = \frac{\partial f(\vec{w}^{(t)})}{\partial \vec{w}} \;^{0} + \frac{\partial \langle \vec{w} - \vec{w}(t), \nabla f(\vec{w}^{(t)}) \rangle}{\partial \vec{w}} + \frac{\partial}{\partial \vec{w}} \frac{\lambda}{2} \| \vec{w} - \vec{w}^{(t)} \|^2$$

$f(\vec{w}^{(t)})$ is not dependent of $\vec{w}$

$$0 = \nabla f(\vec{w}^{(t)})^T + \frac{\lambda}{2} \cdot 2 \cdot (\vec{w} - \vec{w}^{(t)})^T$$

$$\lambda \vec{w} = \lambda w^{(t)} - \nabla f(w^{(t)})$$

$$W = w^{(t)} - \frac{1}{\lambda} \nabla f(w^{(t)}) \quad (A)$$

comparing (1) and (A),

$$\boxed{\eta = \frac{1}{\lambda}} \quad (*)$$

The update rule moves toward the opposite direction of the gradient (page 100 of textbook). And when the gradient is zero, $W$ is optimized.

$\eta$ and $\lambda$ have an inverse relationship. $(*)$.

2|

(eq 5) $\sum_{t=1}^{T} \langle w^{(t)} - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \|v_t\|^2$

$\rightarrow \left[ \sum \left( \langle w^{(t)}, v_t \rangle - \langle w^*, v_t \rangle \right) \right] 2\eta \leq \|w^*\|^2 + \eta^2 \sum_{t=1}^{T} \|v_t\|^2$

$\hookrightarrow = w^{(t+1)} + \eta v_t$ (update rule)

$LHS = \left[ \sum_{t=1}^{T} \langle w^{(t+1)}, v_t \rangle + \eta \langle v_t, v_t \rangle - \langle w^*, v_t \rangle \right] 2\eta$

$= 2\eta \sum_{t=1}^{T} \langle w^{(t+1)} - w^*, v_t \rangle + 2\eta^2 \sum_{t=1}^{T} \|v_t\|^2$

$\leq RHS = \|w^*\|^2 + \boxed{\eta^2 \sum_{t=1}^{T} \|v_t\|^2}$

$\hookrightarrow$ move to LHS
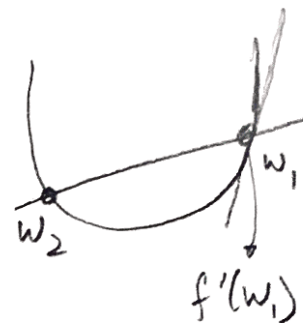
$RHS = \|w^*\|^2 \geq 2\eta \sum_{t=1}^{T} \langle w^{(t+1)} - w^*, v_t \rangle + \eta^2 \sum_{t=1}^{T} \|v_t\|^2 = LHS$

Add $\sum_{t=1}^{T} \|w^{(t+1)} - w^*\|^2$ to both sides

$\|w^*\|^2 + \sum_{t=1}^{T} \|w^{(t+1)} - w^*\|^2 \geq \sum_{t=1}^{T} \|w^{(t+1)} - w^*\|^2 + 2\eta \sum_{t=1}^{T} \langle w^{(t+1)} - w^*, v_t \rangle$

$+ \eta^2 \sum_{t=1}^{T} \|v_t\|^2$

$\|w^*\|^2 + \sum_{t=1}^{T} \|w^{(t+1)} - w^*\|^2 \geq \left( \sum \|w^{(t+1)} - w^*\|^2 + \eta \sum_{t=1}^{T} \|v_t\|^2 \right)^2$

<u>3]</u> $\quad f(w^{(1)}) - f(w^*) \leq \langle w^{(1)} - w^*, f'(w^{(1)}) \rangle$

$\quad f(w^{(2)}) - f(w^*) \leq \langle w^{(2)} - w^*, f'(w^{(2)}) \rangle$

$\quad \vdots$

$\quad f(w^{(T)}) - f(w^*) \leq \langle w^{(T)} - w^*, f'(w^{(T)}) \rangle$



for convex function instantaneous slope is greater than average slope

$$\frac{f(w_1) - f(w_2)}{w_1 - w_2} < f'(w_1)$$

Add all LHS's and RHS's together.

$$\left( \sum_{t=1}^{T} f(w^{(t)}) \right) - T \cdot f(w^*) \leq \sum \langle w^{(t)} - w^*, f'(w^{(t)}) \rangle$$

divide by $T$, $\quad \boxed{\frac{1}{T} \sum_{t=1}^{T} f(w^{(t)})} - f(w^*) \leq \frac{1}{T} \sum_{t=1}^{T} \langle w^{(t)} - w^*, f'(w^{(t)}) \rangle$

by Jensen's Inequality, $\quad f\left( \frac{1}{T} \sum_{t=1}^{T} f(w^{(t)}) \right) \leq \frac{1}{T} \sum_{t=1}^{T} f(w^{(t)})$

So,

$$f\left( \underbrace{\frac{1}{T} \sum_{t=1}^{T} f(w^{(t)})} \right) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^{T} f(w^{(t)}) - f(w^*)$$

$$f(\bar{w}) - f(w^*) \qquad \leq \frac{1}{T} \sum_{t=1}^{T} \langle w^{(t)} - w^*, f'(w^{(t)}) \rangle$$

take result from #2, and divide by $T$.

$$\frac{1}{T} \sum_{t=1}^{T} \langle w^{(t)} - w^*, v_t \rangle \leq \frac{1}{T} \left( \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \|v_t\|^2 \right) \leq \text{---} \rightarrow$$

$$- - - \rightarrow \quad \leq \frac{1}{T}\left(\frac{B^2}{\eta} + \frac{\eta}{2} \cdot T\rho^2\right)$$

Now, rewrite (eq 6) using this inequality

$$f(\bar{w}) - f(w^*) \leq \frac{1}{T}\left(\frac{B^2}{2\eta} + \frac{\eta}{2} \cdot T\rho^2\right)$$

if $\quad \eta = \sqrt{\frac{B^2}{\rho^2 T}} = \frac{B}{\rho}\frac{1}{\sqrt{T}} \quad$ then

$$f(\bar{w}) - f(w^*) \cdot \leq \frac{1}{T}\left(\frac{B^2}{2} \cdot \frac{\rho\sqrt{T}}{B} + \frac{1}{2}\frac{B}{\rho}\frac{1}{\sqrt{T}} \cdot T\rho^2\right)$$

$$\leq \frac{B\rho}{2} \cdot \left(\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{T}}\right)$$

$$\leq \frac{B\rho}{2} \cdot \frac{2}{\sqrt{T}}$$

$$f(\bar{w}) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$$

Note, $\quad \frac{B\rho}{\sqrt{T}} \propto \frac{1}{\sqrt{T}}$

**4)** $f_1(w) = -\ln\left(1 - \frac{1}{1+e^{-w}}\right) = -\ln\left(\frac{e^{-w}}{1+e^{-w}}\right)$

$\frac{d\,f_1(w)}{dw} = (-1) \cdot \left(\frac{1+e^{-w}}{e^{-w}}\right) \cdot \frac{e^{-w}(-1) \cdot (1+e^{-w}) - e^{-w} \cdot e^{-w}(-1)}{(1+e^{-w})^2}$

$= \frac{-(1+e^{-w})}{e^{-w}} \cdot \frac{-e^{-w} - e^{-2w} + e^{-2w}}{(1+e^{-w})^2{}_1}$
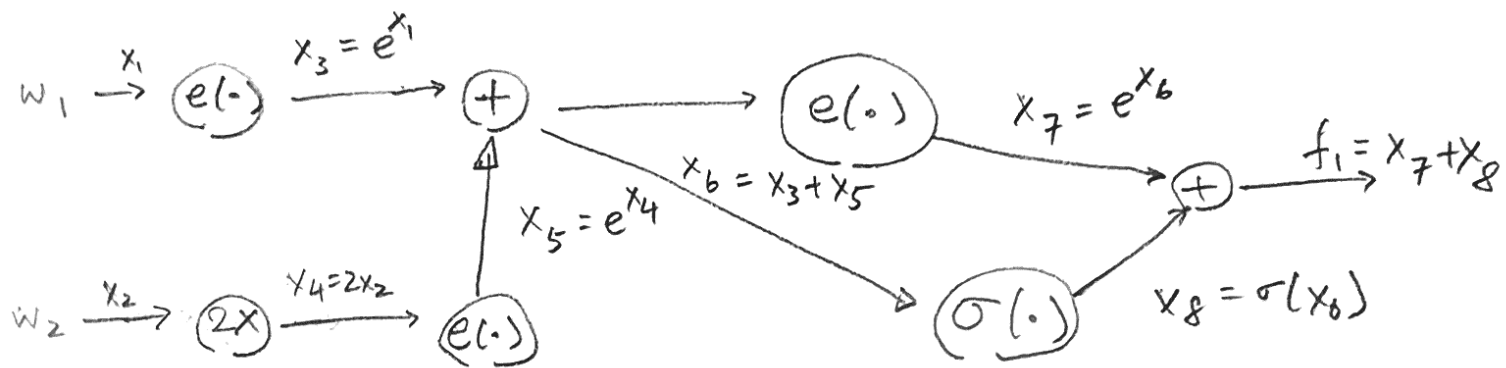
$= \frac{1}{1+e^{-w}} > 0$

$f_2(w) = -\ln\left(\frac{1}{1+e^{-w}}\right)$

$\frac{d\,f_2(w)}{dw} = (-1)(1+e^{-w}) \cdot \frac{(-1)(e^{-w})(-1)}{(1+e^{-w})^2}$

$= \frac{-e^{-w}}{1+e^{-w}} < 0$

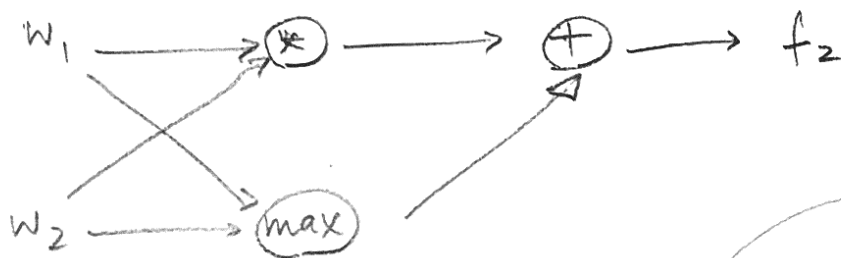There is $\underline{\underline{no}}$ guarantee, since the gradients, depending on

which term is picked, is in opposite directions.

**5]** **a]** $f_1(w_1, w_2) = e^{e^{w_1} + e^{2w_2}} + \sigma\left(e^{w_1} + e^{2w_2}\right)$

$w_1 \xrightarrow{x_1} \boxed{e(\cdot)} \xrightarrow{x_3 = e^{x_1}} \oplus$

$x_5 = e^{x_4}$

$x_6 = x_3 + x_5$

$\boxed{e(\cdot)} \xrightarrow{x_7 = e^{x_6}}$

$\oplus \xrightarrow{f_1 = x_7 + x_8}$

$w_2 \xrightarrow{x_2} \boxed{2x} \xrightarrow{x_4 = 2x_2} \boxed{e(\cdot)}$

$\boxed{\sigma(\cdot)}$  $x_8 = \sigma(x_6)$

$f_1(1, -1) = e^{(e + e^{-2})} + \sigma(e + e^{-2}) = 18.30$

$f_2(w_1, w_2) = w_1 w_2 + \max(w_1, w_2)$

$w_1 \longrightarrow \circledast \longrightarrow \oplus \longrightarrow f_2$

$w_2 \longrightarrow \boxed{\max}$

$f_2(1, -1) = -1 + 1 = 0$

$\therefore f = \begin{bmatrix} 18.30 \\ 0 \end{bmatrix}$

**b]** Jacobian $\Rightarrow J_{i,j} = \dfrac{\partial}{\partial x_j} f_i(x)$

$\begin{bmatrix} \dfrac{\partial f_1}{\partial w_1} & \dfrac{\partial f_1}{\partial w_2} \\[2mm] \dfrac{\partial f_2}{\partial w_1} & \dfrac{\partial f_2}{\partial w_2} \end{bmatrix} = \begin{pmatrix} 48.2 & 4.76 \\[2mm] 0 & 1 \end{pmatrix}$

$*$ these values are computed from the following ---.→

$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

Set $h = \Delta w = 0.01$, $\vec{w} = (1, -1)$

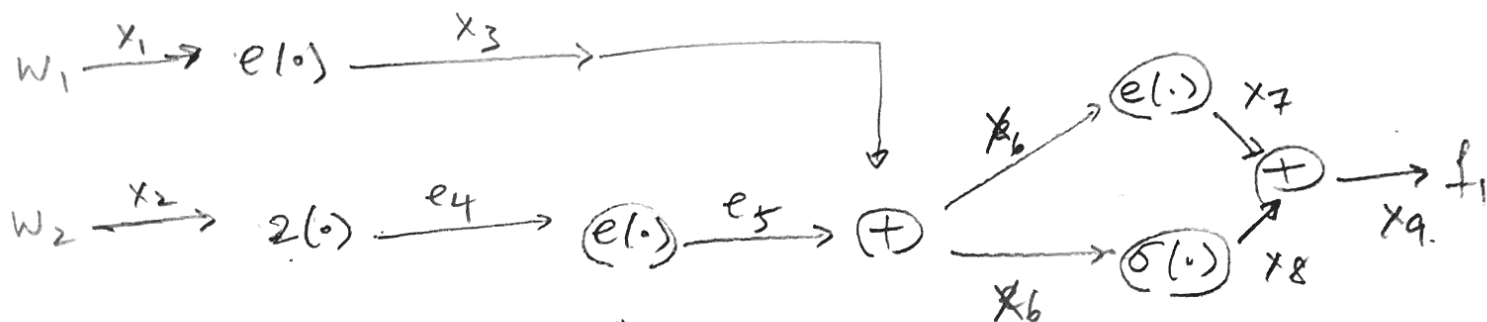$$\frac{\partial f_1}{\partial w_1} = \frac{f_1(1.01, -1) - f_1(1, -1)}{0.01} = 48.192$$

$$\frac{\partial f_1}{\partial w_2} = \frac{f_1(1, -0.99) - f_1(1, -1)}{0.01} = 4.764$$

$$\frac{\partial f_2}{\partial w_1} = \frac{f_2(1.01, -1) - f_2(1, -1)}{0.01} = 0$$

$$\frac{\partial f_2}{\partial w_2} = \frac{f_2(1, -0.99) - f_2(1, -1)}{0.01} = 1$$

**5) c)** $F_1$, forward, so start from $x_1$, $x_2$.

$$W_1 \xrightarrow{x_1} e(\cdot) \xrightarrow{\quad x_3 \quad}$$

$$W_2 \xrightarrow{x_2} 2(\cdot) \xrightarrow{e_4} e(\cdot) \xrightarrow{e_5} (+) \xrightarrow{x_6} e(\cdot) \xrightarrow{x_7} (+) \xrightarrow{x_9} f_1$$

$$\xrightarrow{x_6} \sigma(\cdot) \xrightarrow{x_8} (+)$$

| equations<br><br>Table 1 | $\dfrac{\partial f_1}{\partial W_1}$ | $\dfrac{\partial f_1}{\partial t_2}$ |
|---|---|---|
| $x_1 = W_1$<br>$x_2 = W_2$ | $\mathring{x}_1 = 1$<br>$\mathring{x}_2 = 0$ | $\mathring{x}_1 = 0$<br>$\mathring{x}_2 = 1$ |
| $x_3 = e^{x_1}$ | $\mathring{x}_3 = e^{x_1} \cdot \mathring{x}_1$ | |
| $x_4 = 2x_2$ | $\mathring{x}_4 = 2\mathring{x}_2$ | |
| $x_5 = e^{x_4}$ | $\mathring{x}_5 = e^{x_4} \cdot \mathring{x}_4$ | |
| $x_6 = x_3 + x_5$ | $\mathring{x}_6 = \mathring{x}_3 + \mathring{x}_5$ | |
| $x_7 = e^{x_6}$ | $\mathring{x}_7 = e^{x_6} \cdot \mathring{x}_6$ | |
| $x_8 = \sigma(x_6)$ | $\mathring{x}_8 = \sigma(x_6)\left(1 - \sigma(x_6)\right) \cdot \mathring{x}_6$ | |
| $x_9 = x_7 + x_8$ | $\mathring{x}_9 = \mathring{x}_7 + \mathring{x}_8$ | |

Consider $\dfrac{\partial f_1}{\partial w_1}$, $\dot{x}_1 = 1$, $\dot{x}_2 = 0$ (from chart above ↑)

$$x_1 = w_1 = 1, \quad x_2 = w_2 = -1$$

$$\frac{\partial f_1}{\partial w_1} = \dot{x}_9 = \dot{x}_7 + \dot{x}_8 = \left(e^{x_6} \cdot \dot{x}_6\right) + \left[\sigma(x_6)\left(1 - \sigma(x_6)\right) \dot{x}_6\right]$$

$$\dot{x}_6 = \dot{x}_3 + \dot{x}_5 = \left(e^{x_1} \cdot \dot{x}_1\right) + \left(e^{x_4} \cdot \dot{x}_4\right)$$

$$= \left(e^{x_1} \cdot \dot{x}_1\right) + \left(e^{x_4} \cdot 2\dot{x}_2\right)$$

$$= \left(e^{x_1} \cdot 1\right) + \left(e^{x_4} \cdot 2(0)\right) = e^{x_1}$$

So, $\dfrac{\partial f_1}{\partial w_1} = e^{x_6}\left(e^{x_1}\right) + \left[\sigma(x_6)(1 - \sigma(x_6))\right]e^{x_1}$

$$x_1 = w_1$$

($\maltese\maltese$) $\quad x_6 = x_3 + x_5 = e^{x_1} + e^{x_4} = e^{x_1} + e^{2x_2} = e^{w_1} + e^{2w_2}$

(A) So, $\dfrac{\partial f_1}{\partial w_1} = e^{\left(e^{w_1} + e^{2w_2}\right)} \cdot e^{w_1} + \sigma\left(e^{w_1} + e^{2w_2}\right)\left(1 - \sigma\left(e^{w_1} + e^{2w_2}\right)\right)e^{w}$

Consider $\dfrac{\partial f_1}{\partial w_2}$, $\dot{x}_1 = 0$, $\dot{x}_2 = 1$

$$\dot{x}_3 = 0, \quad \dot{x}_4 = 2, \quad \dot{x}_5 = 2e^{x_4}, \quad \dot{x}_6 = \dot{x}_5, \quad \dot{x}_7 = e^{x_6} \cdot 2e^{x_4}$$

$$\dot{x}_8 = \sigma(x_6)\left(1 - \sigma(x_6)\right) \cdot 2e^{x_4}$$

$$\dot{x}_9 = e^{x_6} \cdot 2e^{x_4} + \sigma(x_6)\left(1 - \sigma(x_6)\right) \cdot 2e^{x_4}$$

(B) $\therefore \dfrac{\partial f_1}{\partial w_2} = e^{\left(e^{w_1} + e^{2w_2}\right)} \cdot 2e^{2w_2} + \sigma\left(e^{w_1} + e^{2w_2}\right)\left(1 - \sigma\left(e^{w_1} + e^{2w_2}\right)\right)2e^{2w_2}$

**5] c] continued.**   F2, forward



$$\frac{\partial f_1}{\partial w_1} = 47.3 \quad \text{from (A),}$$

$$\frac{\partial f_1}{\partial w_2} = 4.71 \quad \text{from (B),}$$

$$\frac{\partial f_2}{\partial w_1} = 0 \quad \text{from (C),}$$

$$\frac{\partial f_2}{\partial w_2} = 1 \quad \text{from (D),}$$

So, Jacobian is:

$$\begin{pmatrix} 47.3 & 4.71 \\ 0 & 1 \end{pmatrix}$$

**Table 2**

| equation | $\dfrac{\partial f w_1}{\partial w_1}$ | $\dfrac{\partial f_2}{\partial w_2}$ |
|---|---|---|
| $x_1 = w_1$ | $\dot{x}_1 = 1$ | $\dot{x}_1 = 0$ |
| $x_2 = w_2$ | $\dot{x}_2 = 0$ | $\dot{x}_2 = 1$ |
| $x_3 = x_1 \cdot x_2$ | $\dot{x}_3 = \dot{x}_1 x_2 + x_1 \dot{x}_2$ | |
| $x_4 = \max(x_1, x_2)$ | $\dot{x}_4 = \begin{cases} \dot{x}_1 & \text{if } x_1 \geq x_2 \\ \dot{x}_2 & \text{o.w.} \end{cases}$ | |
| $x_5 = x_3 + x_4$ | $\dot{x}_5 = \dot{x}_3 + \dot{x}_4$ | |

Consider $\dfrac{\partial f_2}{\partial w_1}$, $\dot{x}_1 = 1$, $\dot{x}_2 = 0$

$$\dot{x}_3 = x_2, \quad \dot{x}_4 = \dot{x}_1 = 1$$

$$\frac{\partial f_2}{\partial w_1} = \dot{x}_5 = x_2 + 1 = w_2 + 1 \quad (c)$$

Consider $\dfrac{\partial f_2}{\partial w_2}$, $\dot{x}_1 = 0$, $\dot{x}_2 = 1$

$$\dot{x}_3 = x_1 = w_1$$

$$\dot{y}_4 = \dot{x}_1 = 0$$

$$\frac{\partial f_2}{\partial w_2} = \dot{x}_5 = w_1 + 0 = w_1 \quad (D)$$

## 5] d] F1, backward (so start at $\dot{x}_9$).

$\dot{x}_9 = 1, \qquad x_9 = x_7 \oplus x_8$, so pass gradient to both.

$\therefore \dot{x}_8 = \dot{x}_9 = 1$ and $\dot{x}_7 = \dot{x}_9 = 1$

Since $x_6$ goes into 2 nodes, $\dot{x}_6$ must sum both gradients.

$\therefore \dot{x}_6 = e^{x_6} \cdot \dot{x}_7 + \sigma(x_6)(1 - \sigma(x_6))\dot{x}_8$

$\qquad = e^{x_6} + \sigma(x_6)(1 - \sigma(x_6))$

$\dot{x}_5 = \dot{x}_6$, and $\dot{x}_3 = \dot{x}_6$, since $x_6 = x_3 \oplus x_5$

$\dot{x}_4 = e^{x_4} \cdot \dot{x}_5 = e^{x_4}\dot{x}_6 = e^{x_4}\left(e^{x_6} + \sigma(x_6)(1 - \sigma(x_6))\right)$

$\dot{x}_2 = 2\dot{x}_4 = 2e^{x_4}\left(e^{x_6} + \sigma(x_6)(1 - \sigma(x_6))\right)$

$\dot{x}_1 = e^{x_1} \cdot \dot{x}_3 = e^{x_1} \cdot \dot{x}_6 = e^{x_1}\left(e^{x_6} + \sigma(x_6)(1 - \sigma(x_6))\right)$

(a) $\dfrac{\partial f_1}{\partial w_1} = \dot{x}_1 = e^{w_1}\left(e^{(e^{w_1} + e^{2w_2})} + \sigma\left(e^{w_1} + e^{2w_2}\right)\left(1 - \sigma(e^{w_1} + e^{2w_2})\right)\right)$

> from (**) of part c, we know $x_6$.

(b) $\dfrac{\partial f_1}{\partial w_2} = \dot{x}_2 = 2e^{2w_2}\left(e^{(e^{w_1} + e^{2w_2})} + \sigma\left(e^{w_1} + e^{2w_2}\right)\left(1 - \sigma(e^{w_1} + e^{2w_2})\right)\right)$

5] d] continued.   Backward, $F_2$.

$\overset{\circ}{x}_5 = 1$ ,   $\dot{x}_4 = \dot{x}_5 = 1$,   $\dot{x}_3 = \dot{x}_5 = 1$

$\dot{x}_1 = x_2 \dot{x}_3 + \overset{\circ}{x}_4 = x_2 = w_2 + 1$

$\overset{\circ}{x}_2 = x_1 \dot{x}_3 + 0 = w_1$

(c) $\dfrac{\partial f_2}{\partial w_1} = w_2 + 1 = \overset{\circ}{x}_1$ ,   (d) $\dfrac{\partial f_2}{\partial w_2} = w_1 = \overset{\circ}{x}_2$.

Note   $a, b, c, d = A, B, C, D$
of part d   of part c.

i.e.   Forward  &  Backward produce same formula for each
entry in ~~Glassian~~ Jacobian

So,   $J = \begin{pmatrix} 47.3 & 4.71 \\ 0 & 1 \end{pmatrix}$   "same as part c"

5] e]   Yes !   I do!