# Autograd

# What AutoDiff Is Not

* Finite differences

$\rightarrow$ One sided:

$$\frac{\partial}{\partial x_i} f(x_1, \ldots, x_N) \approx \frac{f(x_1, \ldots x_i+h, \ldots) - f(x_1, \ldots x_i, x_N)}{h}$$

$\rightarrow$ Or Two sided

$$\frac{\partial}{\partial x_i} f(x_1, \ldots, x_N) \approx \frac{f(x_1, \ldots x_i+h, \ldots) - f(x_1, \ldots, x_i-h, \ldots)}{2h}$$

* Challenges with finite differences

  → Expensive: compute forward pass for each variable

  → Numerically unstable

# Computational Graphs

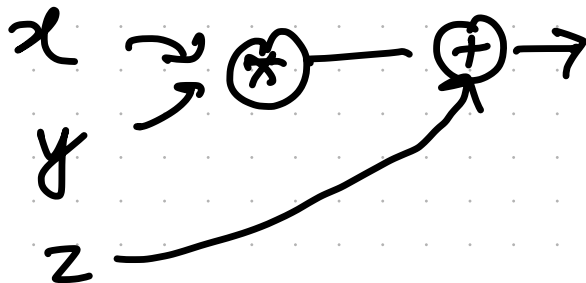* Nodes : operations $(+, *, \dots)$

* Edges : variables / Tensors

# Computational Graphs
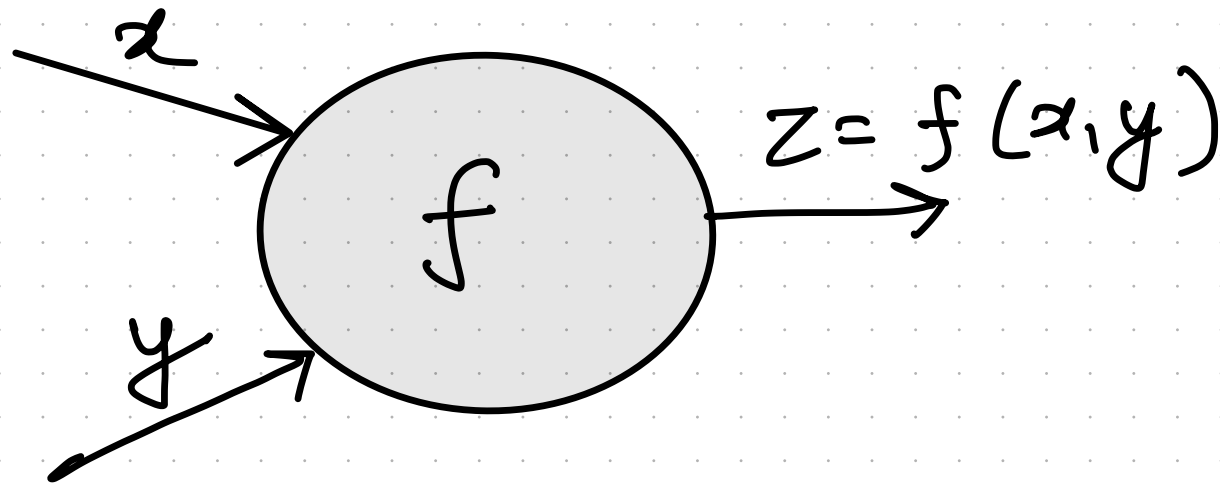
* Nodes : operations $(+, *, \ldots)$

* Edges : variables / Tensors
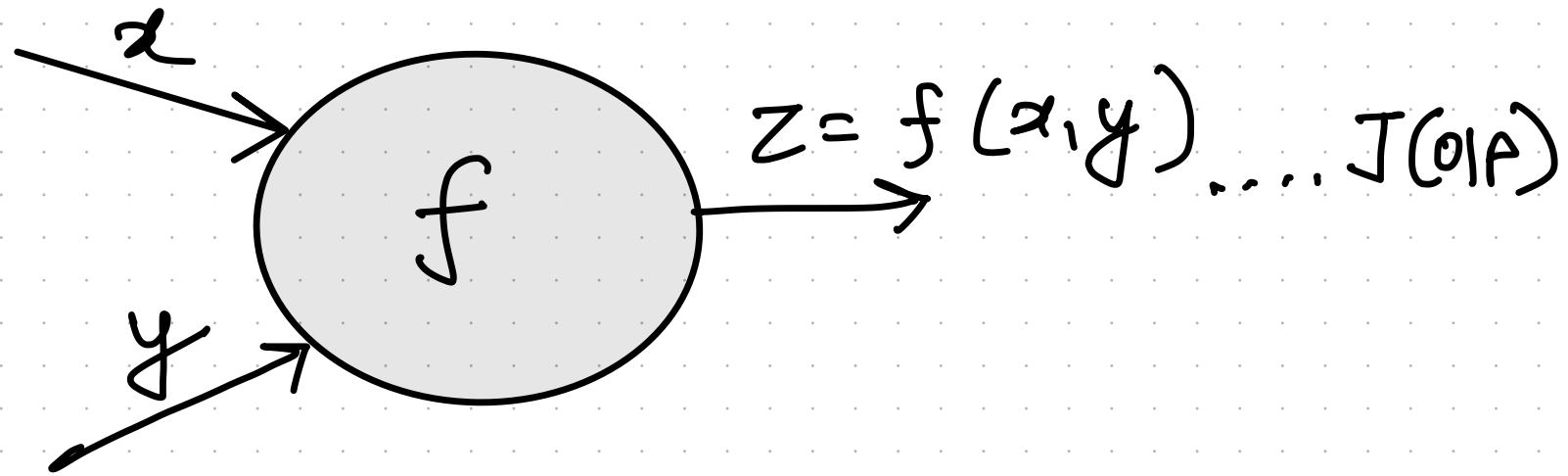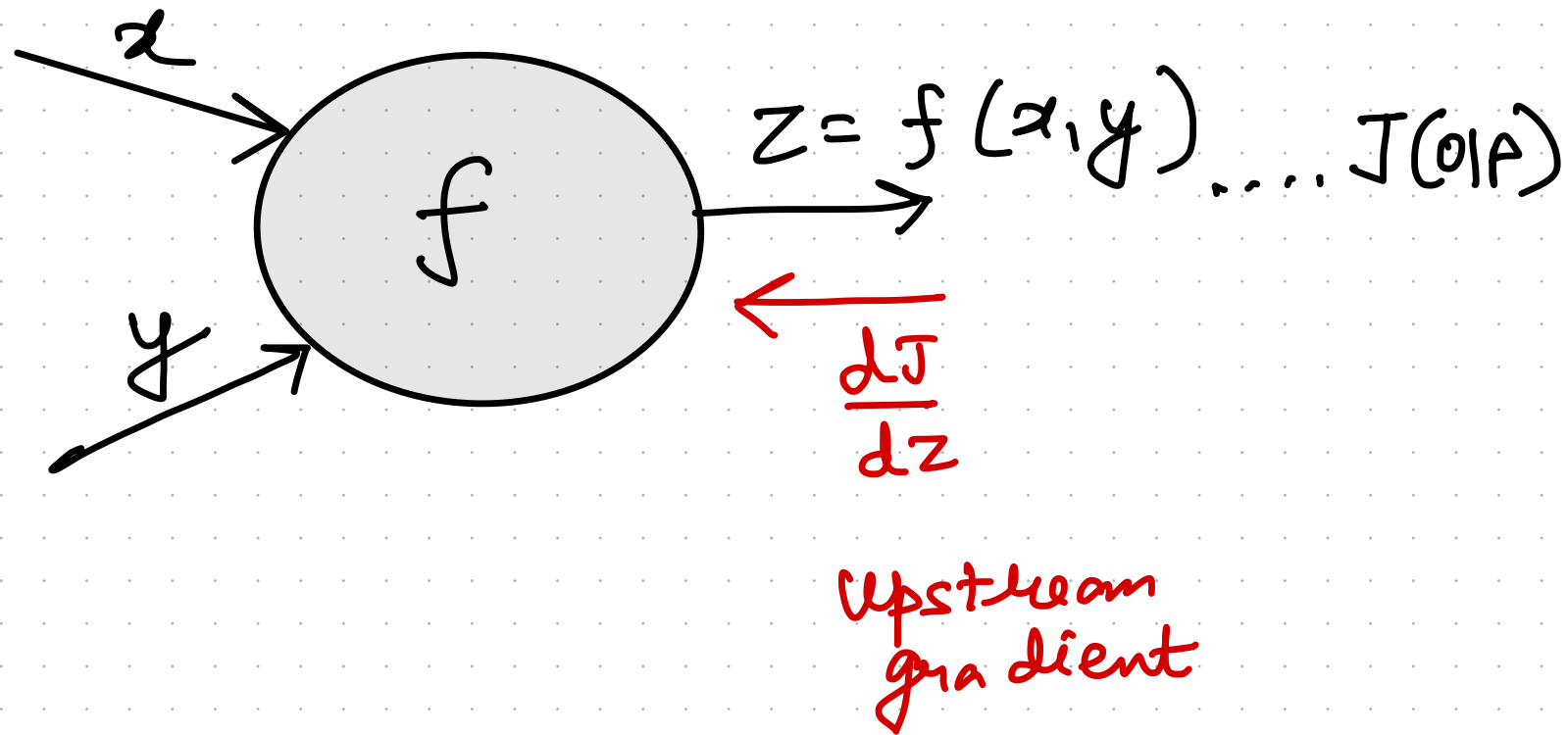                  (and data dependencies)

Example : $(x * y) + z$
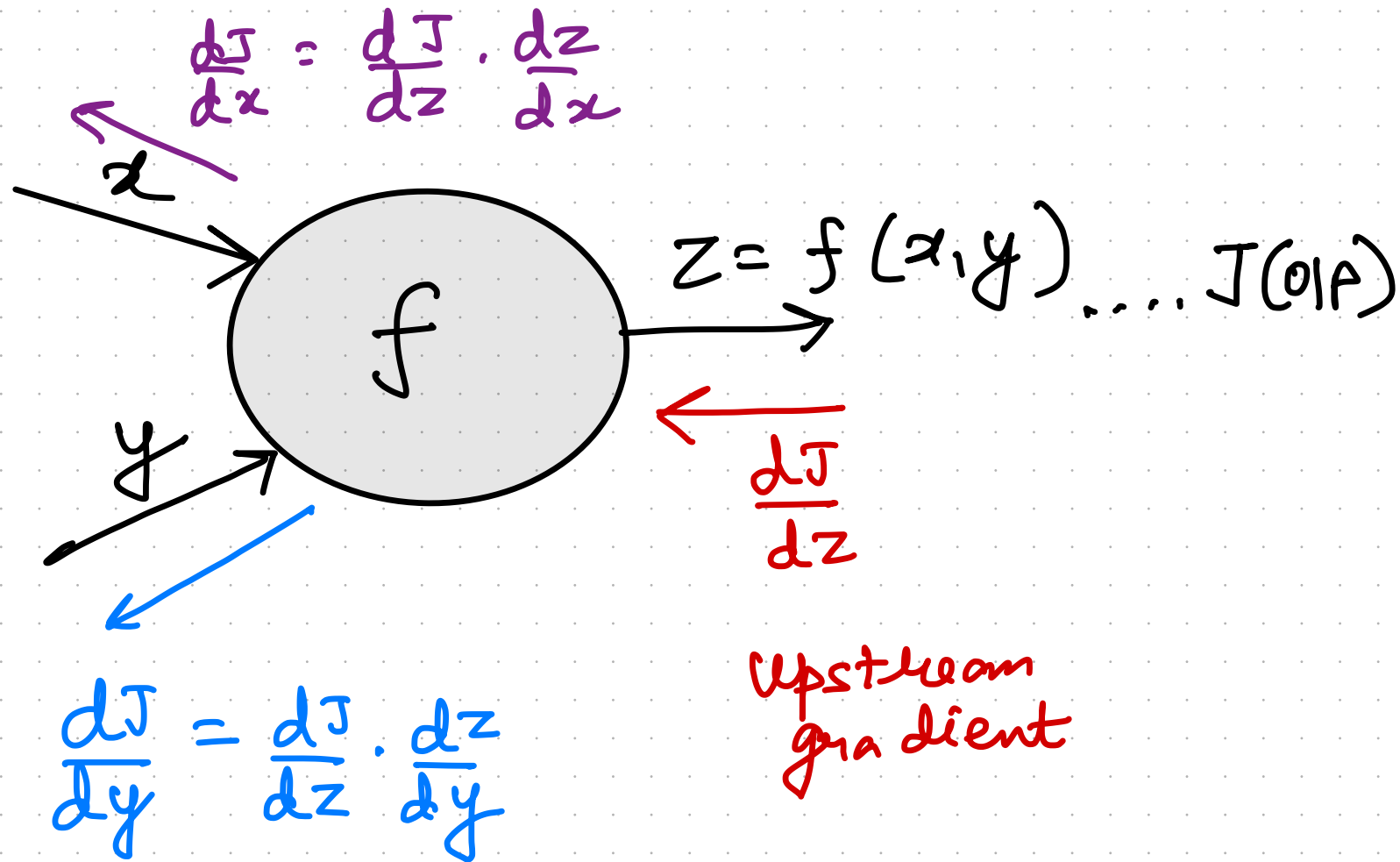
# Back Prop Through Computational Graph

$$z = f(x, y)$$

# Back Prop Through Computational Graph



$$z = f(x, y) \ldots J(O/P)$$

# Back Prop Through Computational Graph



$$z = f(x, y) \ldots J(olp)$$

$$\frac{dJ}{dz}$$

Upstream gradient

# Back Prop Through Computational Graph

$$\frac{dJ}{dx} = \frac{dJ}{dz} \cdot \frac{dz}{dx}$$

$x$

$$z = f(x,y) \ldots\ldots J(O/P)$$

$y$

$$\frac{dJ}{dz}$$

Upstream gradient

$$\frac{dJ}{dy} = \frac{dJ}{dz} \cdot \frac{dz}{dy}$$

# Back Prop Through Computational Graph

$$\frac{dJ}{dx} = \frac{dJ}{dz} \cdot \frac{dz}{dx} \quad \leftarrow \text{Local grad.}$$

$x$

$\leftarrow$ upstream grad.

$$Z = f(x,y) \ldots J(\text{O/P})$$

$y$

$f$

$$\frac{dJ}{dz}$$

Upstream gradient

$$\frac{dJ}{dy} = \frac{dJ}{dz} \cdot \frac{dz}{dy}$$

upstream gradient

Local grad.

# Back Prop Through Computational Graph

$$\frac{dJ}{dx} = \frac{dJ}{dz} \cdot \frac{dz}{dx}$$

$$x$$

$$f$$

$$z = f(x,y) \ldots J(O|P)$$

$$\frac{dJ}{dz}$$

Upstream gradient

$$y$$

$$\frac{dJ}{dy} = \frac{dJ}{dz} \cdot \frac{dz}{dy}$$
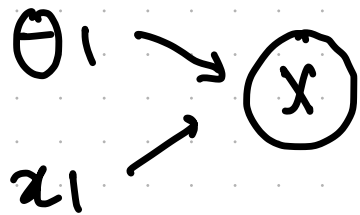
DOWNSTREAM GRADIENT
= UPSTREAM GRADIENT * LOCAL GRADIENT

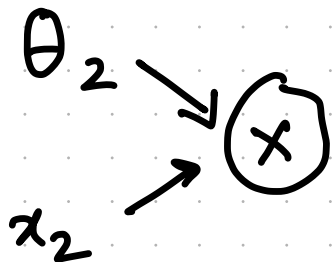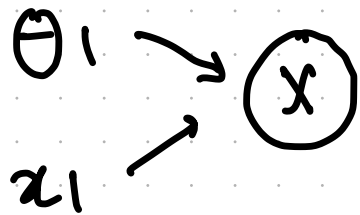$$\hat{y} = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}$$

$$y = 1$$

$$\text{Loss} = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

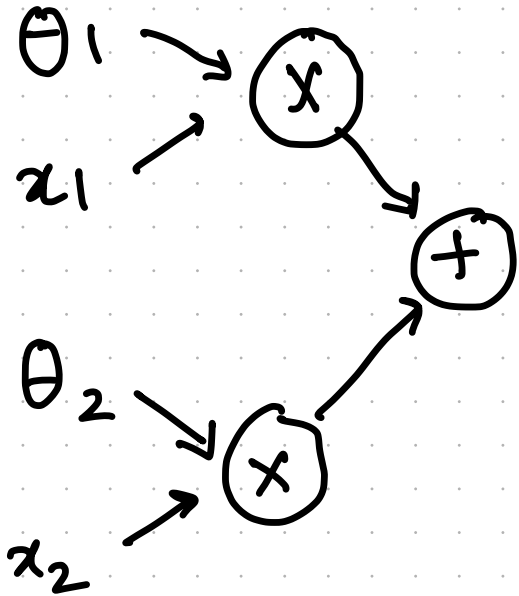$$= -\log \left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$\theta_1$

$x_1$

$\bigotimes$ (X)

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$\theta_1$

$x_1$

$\times$

$\theta_2$

$x_2$

$\times$

$$\text{LOSS} = -1*\log\left(\frac{1}{1+e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$\theta_1$

$x_1$

$\theta_2$

$x_2$

( × ) ( + ) ( × )

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$
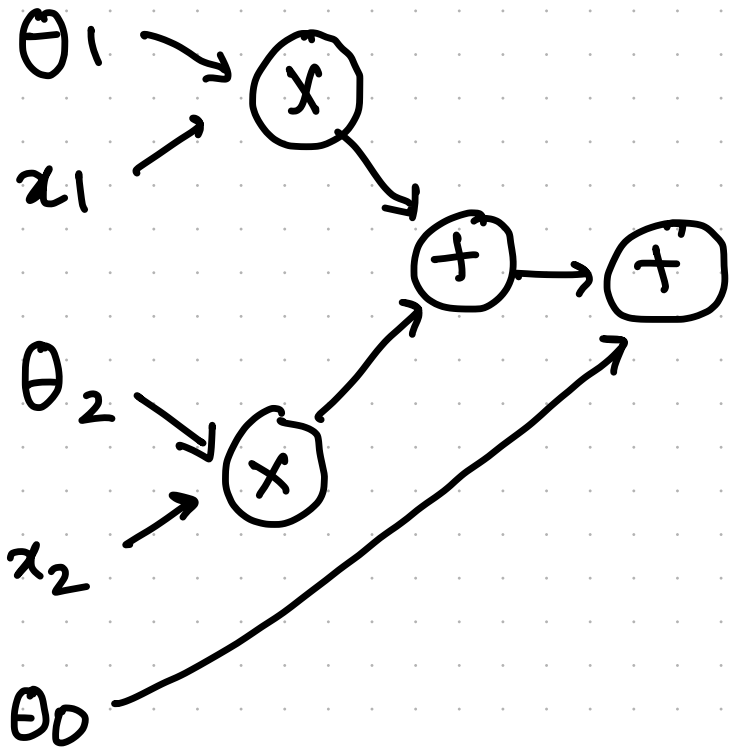
$\theta_1$

$x_1$

$\theta_2$

$x_2$

$\theta_0$

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

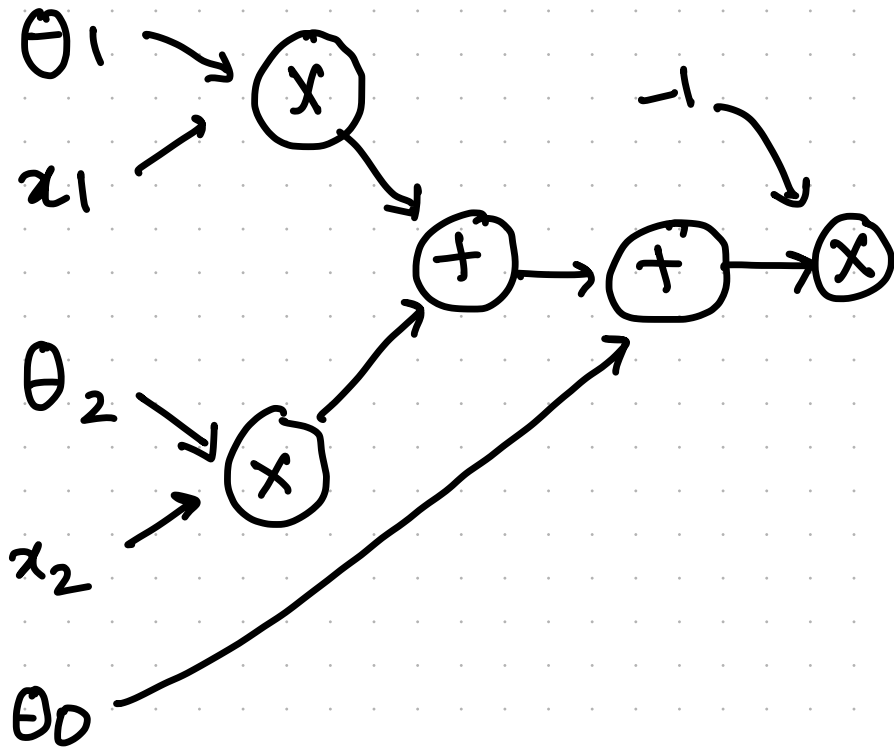$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$
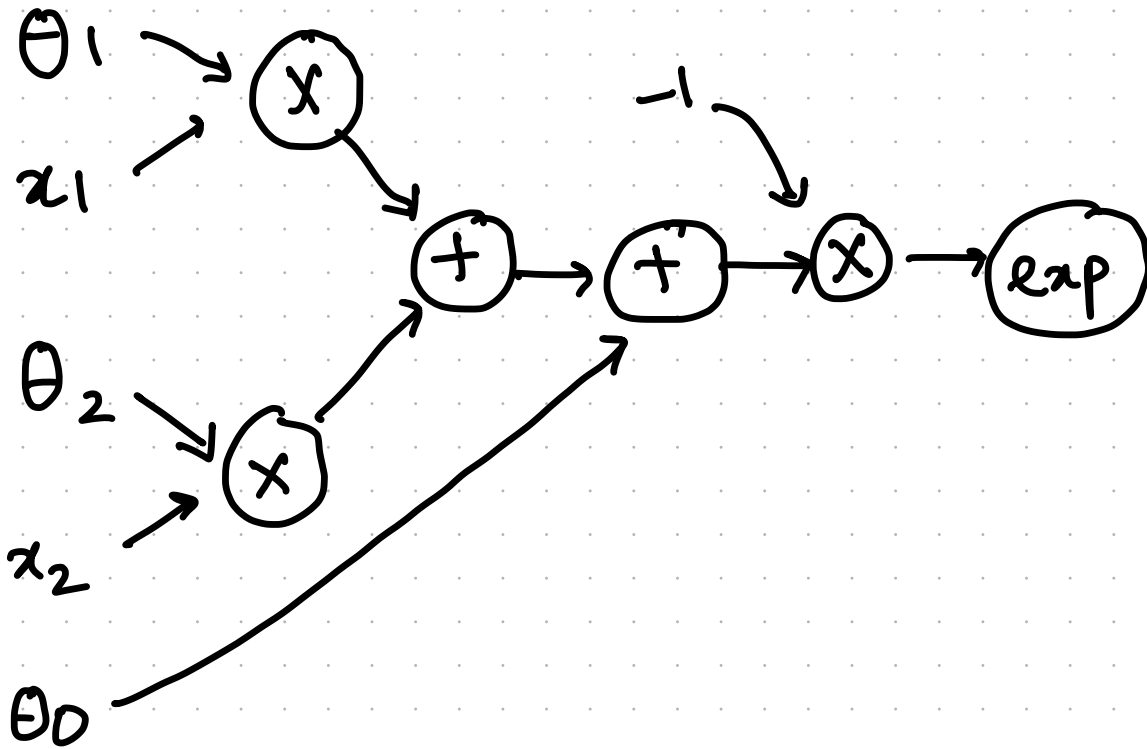
$\theta_1$

$x_1$

$\theta_2$

$x_2$

$\theta_0$

-1

×

×

+

+

×

exp

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$\theta_1$

$x_1$

$\theta_2$

$x_2$

$\theta_0$

$-1$

$1$

$\times$ $+$ $+$ $\times$ $\exp$ $+$ $\times$

$$LOSS = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

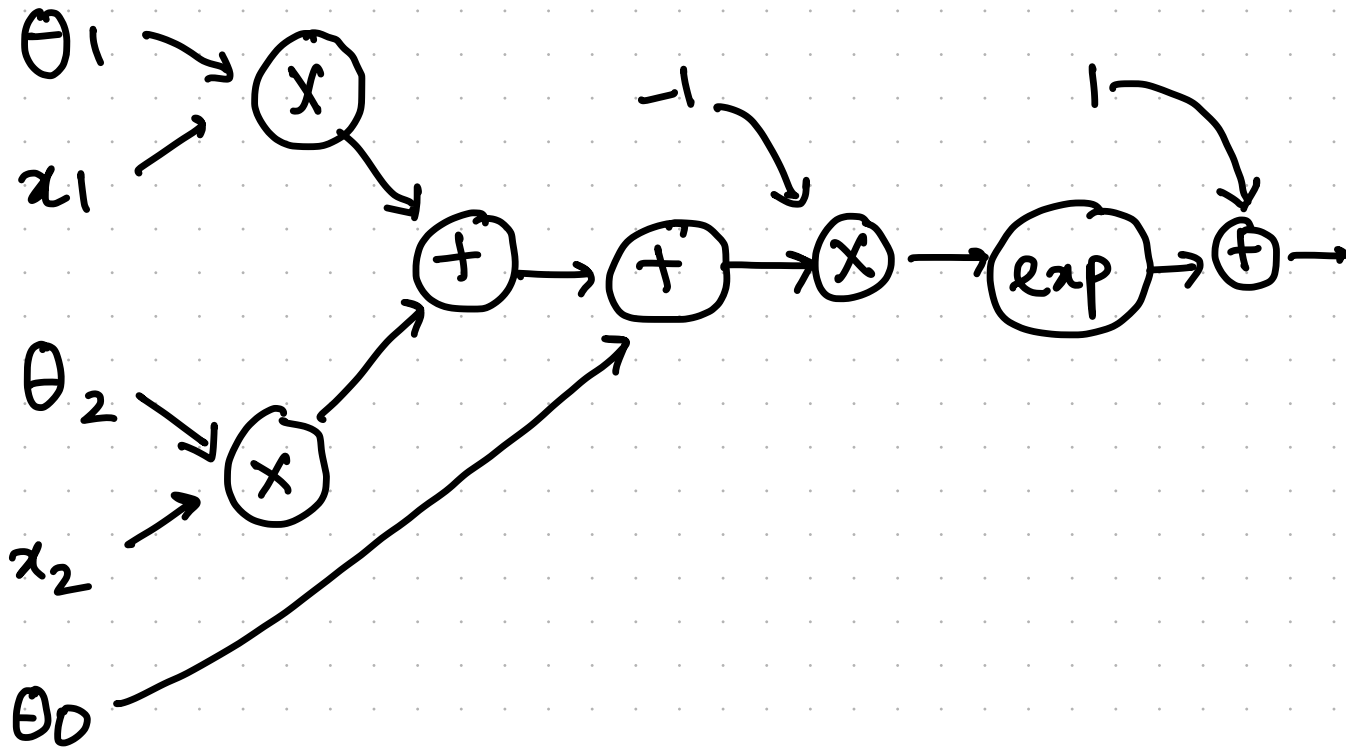$$LOSS = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$LOSS = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$
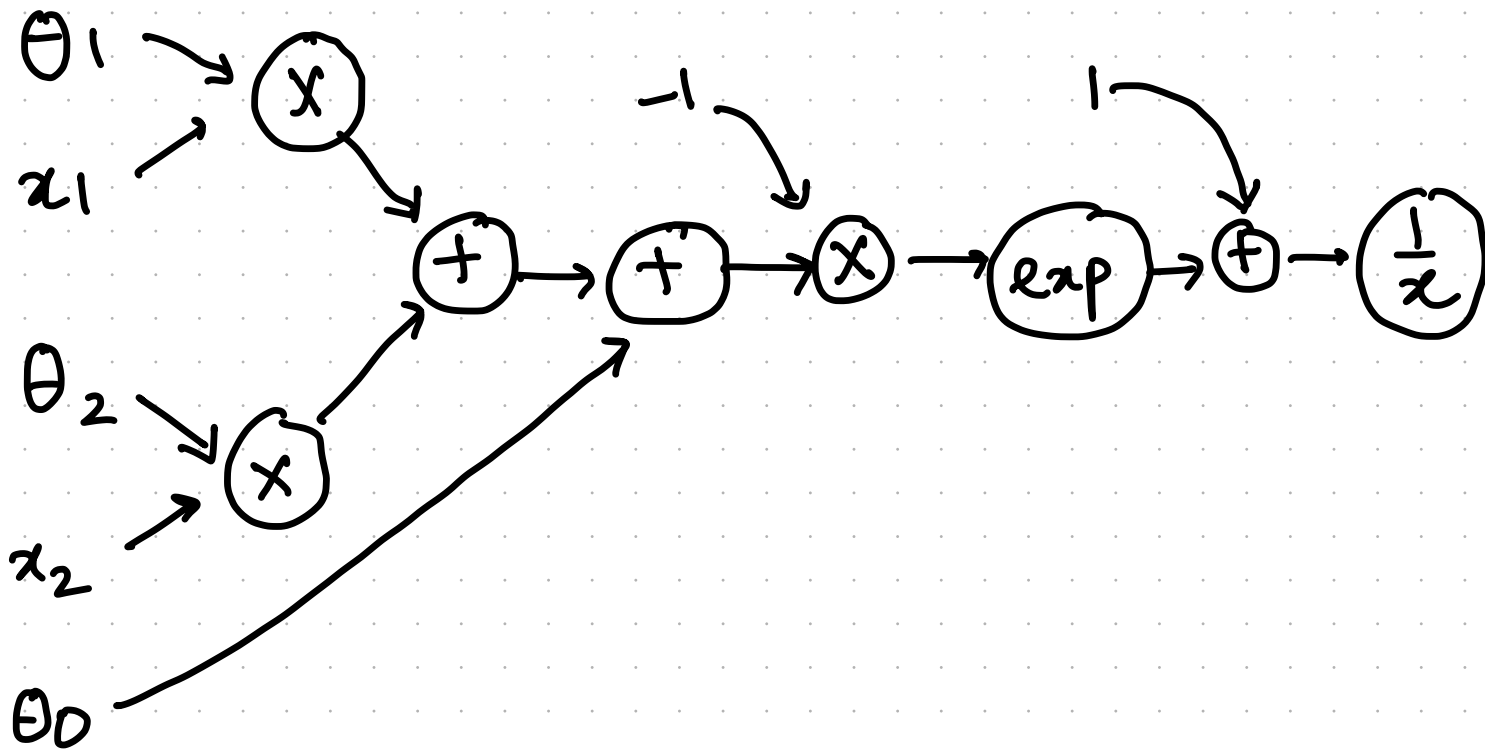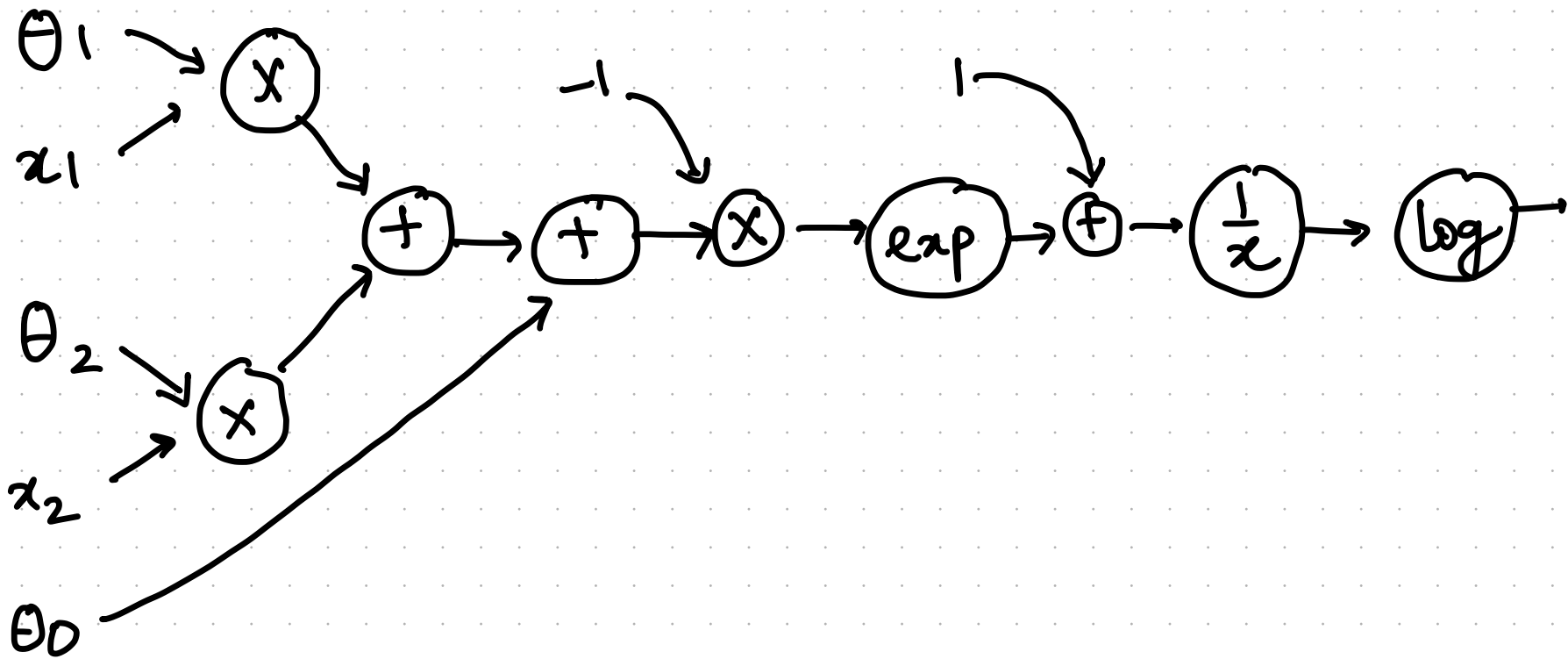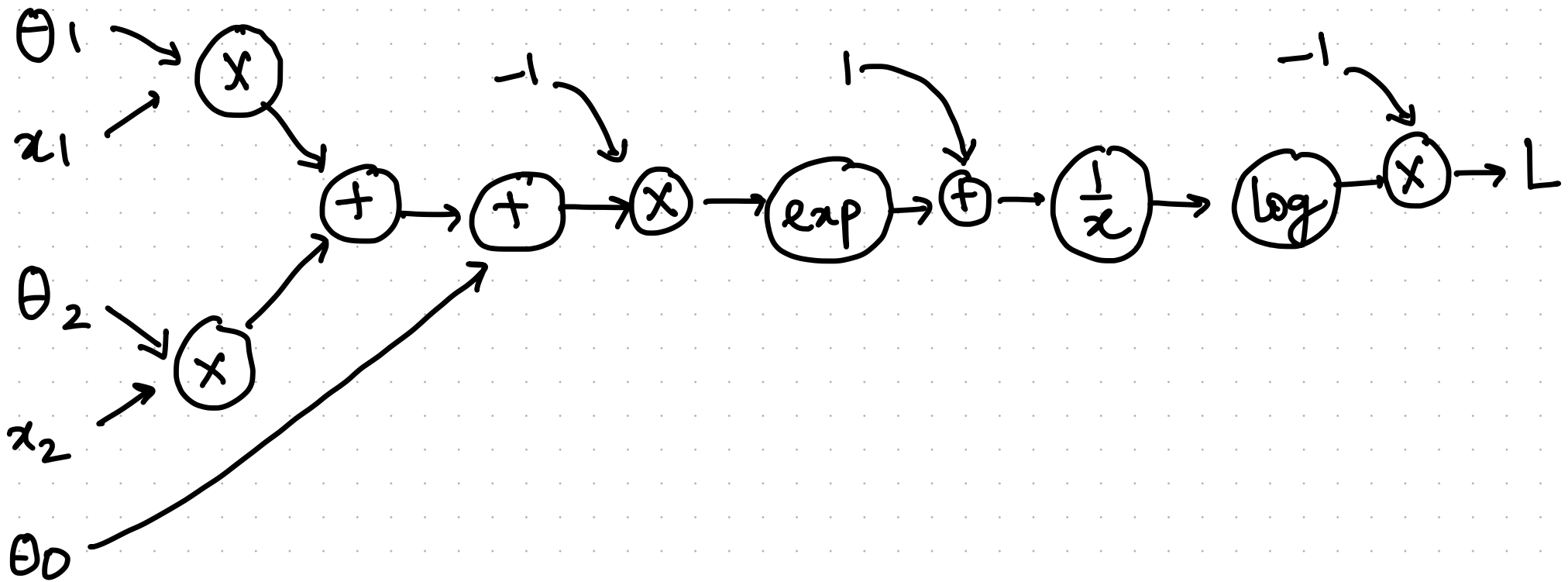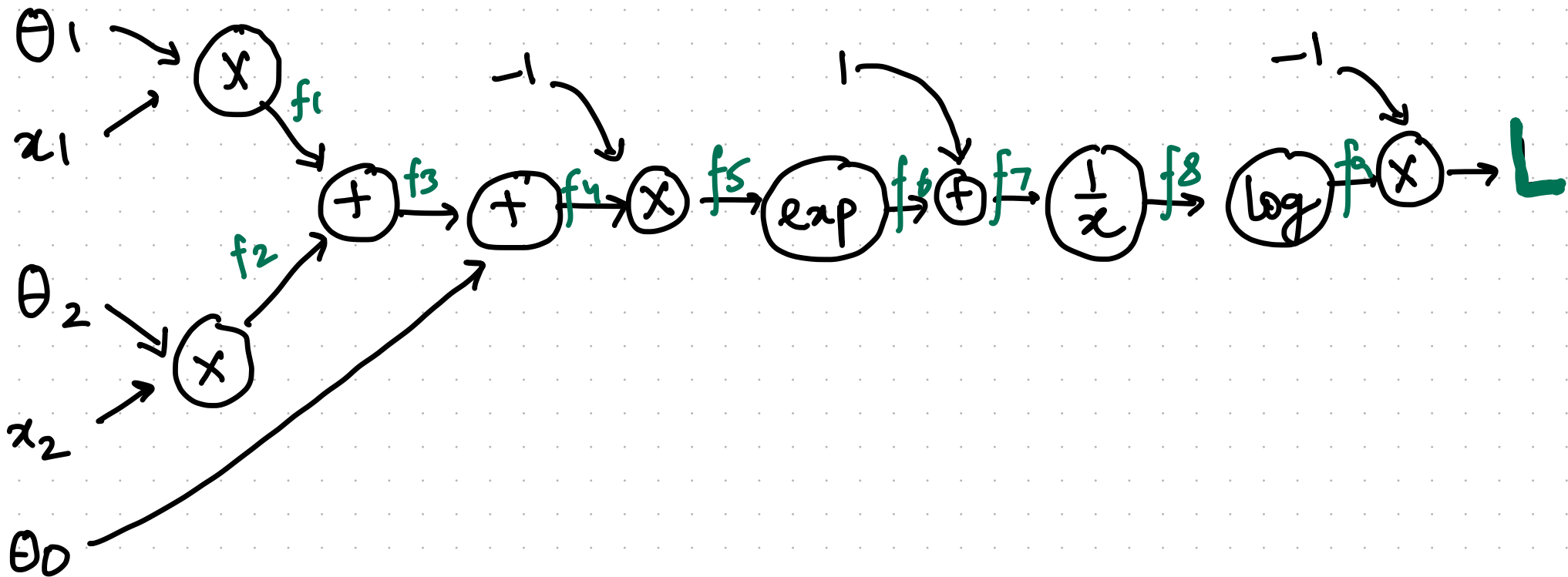
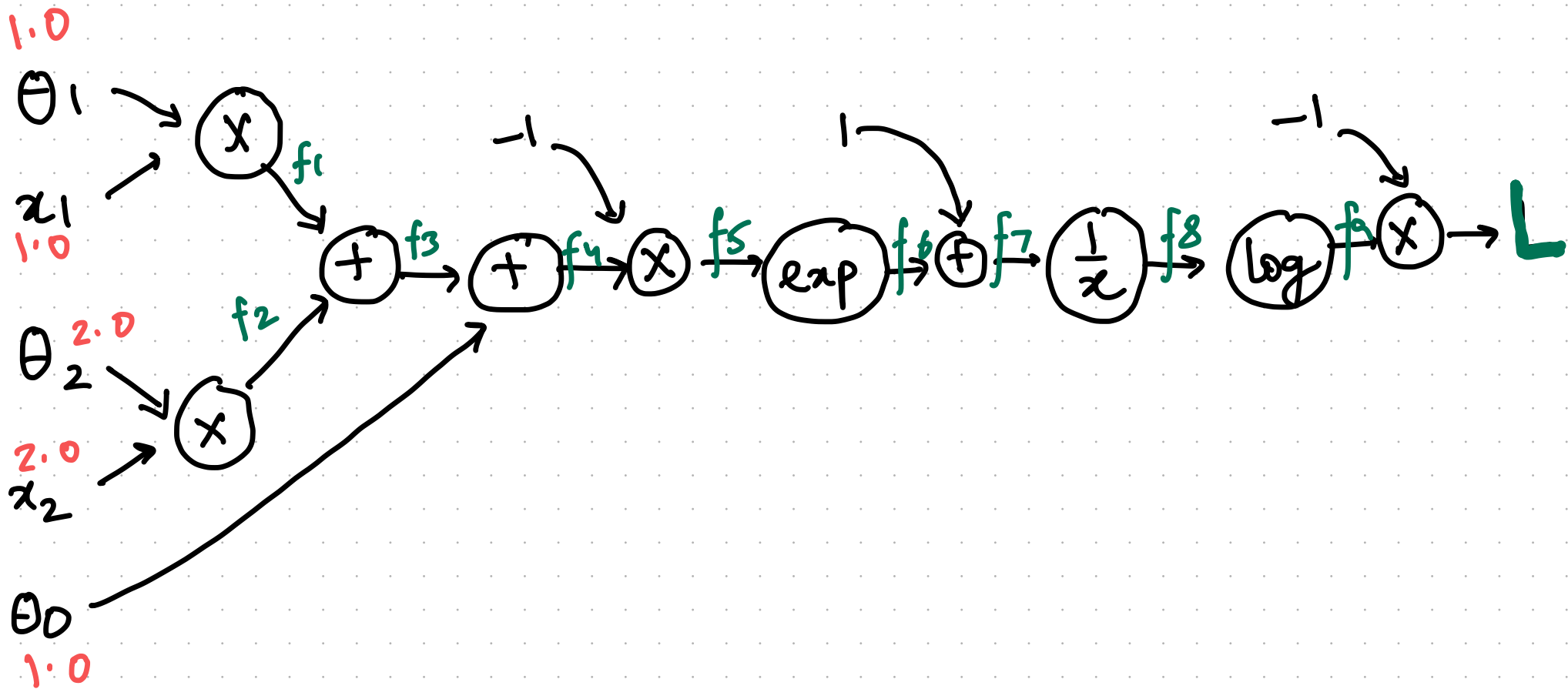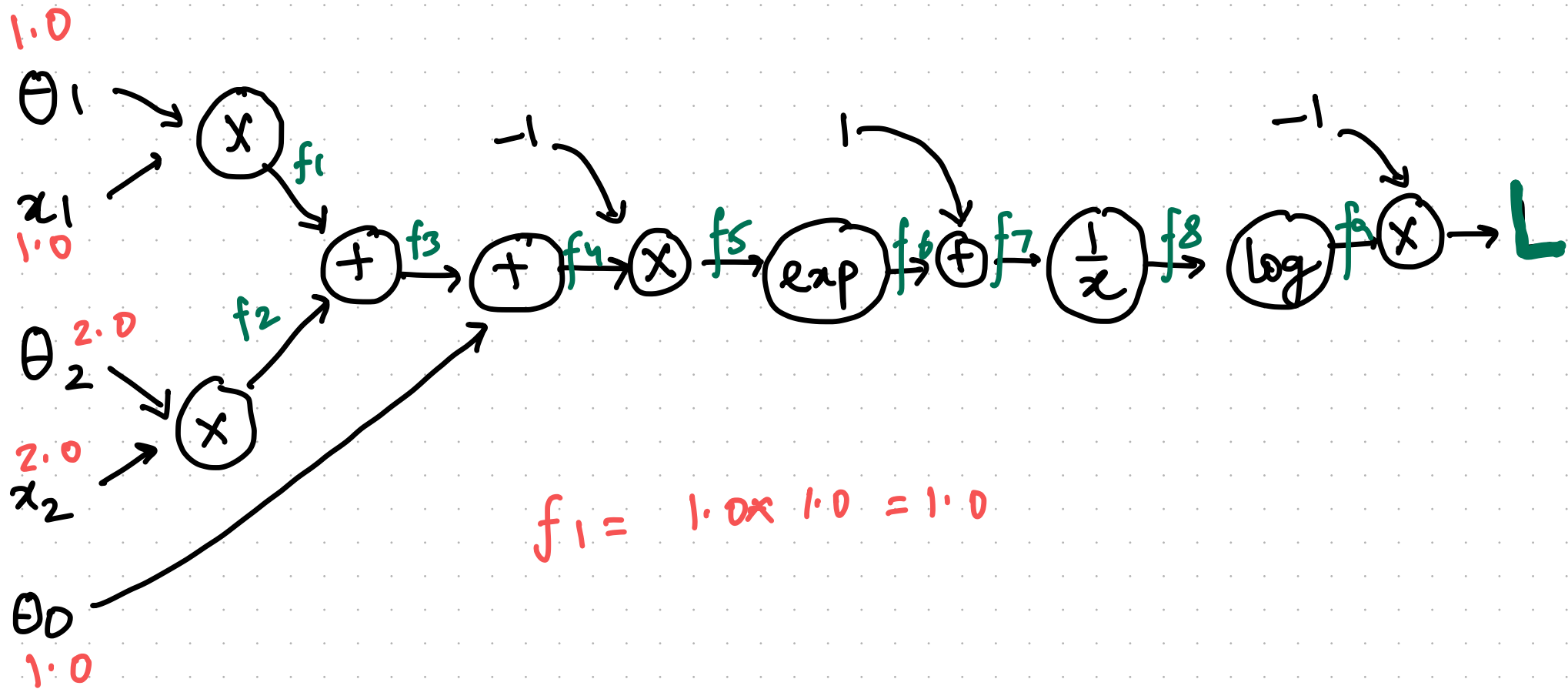$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

1.0
$\theta_1$

$x_1$
1.0

2.0
$\theta_2$

2.0
$x_2$

$\theta_0$
1.0

$f_1 = 1.0 \times 1.0 = 1.0$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$



1.0
$\theta_1$

$x_1$
1.0

$\theta_2$ 2.0

2.0
$x_2$

$\theta_0$

1.0

$f_1 = 1.0$

$f_2 = 4.0$

$f_3 = 5.0$

$-1$

$f_4$  6.0

$f_5$  $-6.0$

exp

$f_6$  .0024

$= 1.0024$

$f_7$

$\frac{1}{x}$  .99

$= 1.0024$

$f_8$  .99

log

$f_9$  $= -0.02$

$-1$

$\times$

$L$

$-0.02$

$= .00247$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$



1.0
$\theta_1$

$x_1$
1.0

$\theta_2$ 2.0

2.0
$x_2$

$\theta_0$

1.0

$f_1 = 1.0$

$-1$

$f_3 = 5.0$

$f_2 = 4.0$

$f_4$

$f_5$

$6.0$

$-6.0$

exp

$f_6$

.00
24

1

$\frac{1}{x}$

= 1.0024

$f_7$

$f_8$

.00
24

.99

log

$-1$

$f_9$

= -0.02

L

=
.00247

$\frac{\partial L}{\partial 2} = 1$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$



1.0
$\theta_1$

$x_1$
1.0

$\theta_2$ 2.0

2.0
$x_2$

$\theta_0$
1.0

$f_1 = 1.0$

$-1$

$f_3 = 5.0$

$f_2 = 4.0$

$f_4$
6.0

$f_5$
$-6.0$

exp

$f_6$
.00 24

$= 1.0024$

$f_7$

$\frac{1}{x}$
.99

$f_8$

log
$= -0.02$

$-1$

$f_9$

$\times$

L
$= .00247$

$$\frac{\partial L}{\partial 2} = 1$$

$$LOSS = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$-1$$

$$f9 \rightarrow \boxed{\times} \rightarrow L$$

$$\xleftarrow{2} \qquad \xleftarrow{1.0}$$

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

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$-1$

$f9 \xrightarrow{} \bigotimes \xrightarrow{} L$

$\xleftarrow{2} \quad \xleftarrow{1.0}$

$= .00247$
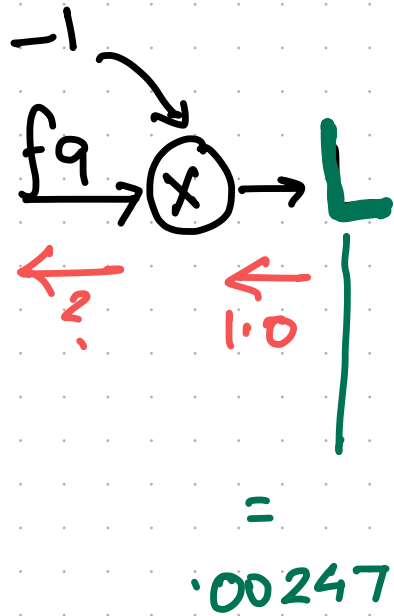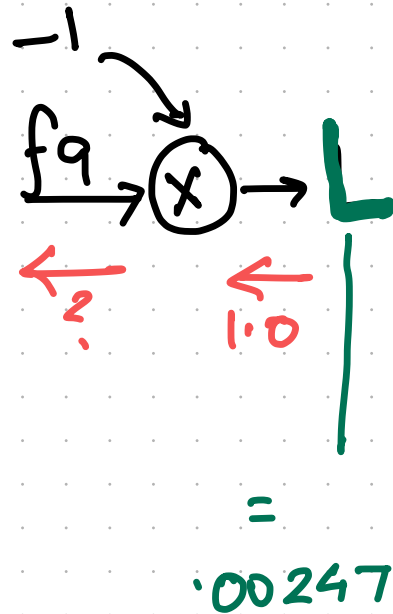
$\frac{\partial L}{\partial L} = 1$

upstream gradient $= 1.0$

$$LOSS = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$
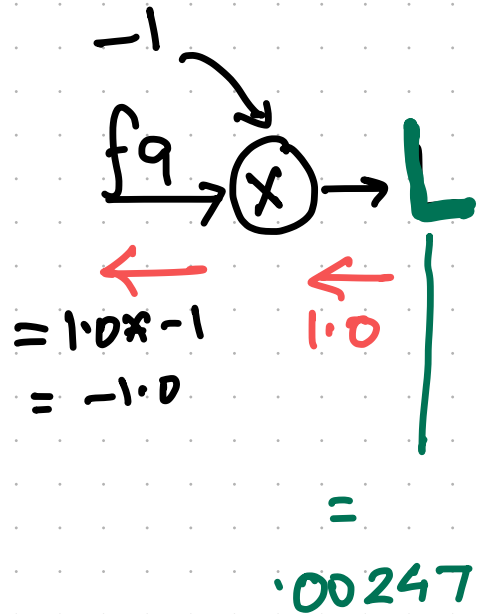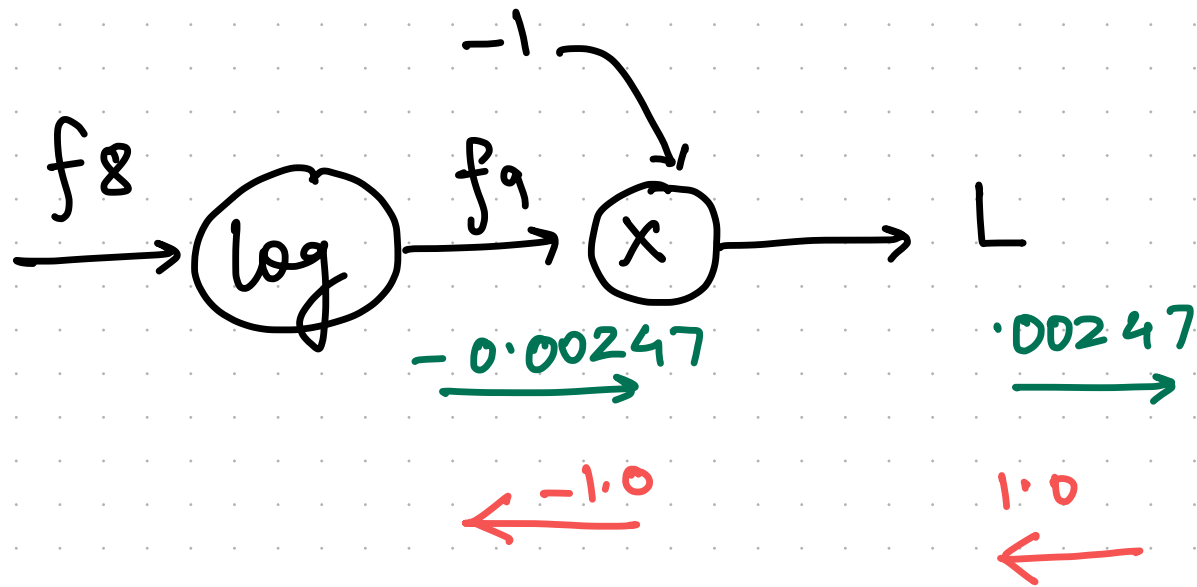


$$\frac{\partial L}{\partial L} = 1$$

upstream gradient = 1.0

$$= .00247$$

$$L = f_9 * -1$$

$$\frac{\partial L}{\partial f_9} = -1 \qquad LOCAL \ GRADIENT = -1$$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$-1$$

$$fg \rightarrow \otimes \rightarrow L$$

$$= 1.0 * -1$$
$$= -1.0$$

$$1.0$$

$$= .00247$$

$$\frac{\partial L}{\partial L} = 1$$
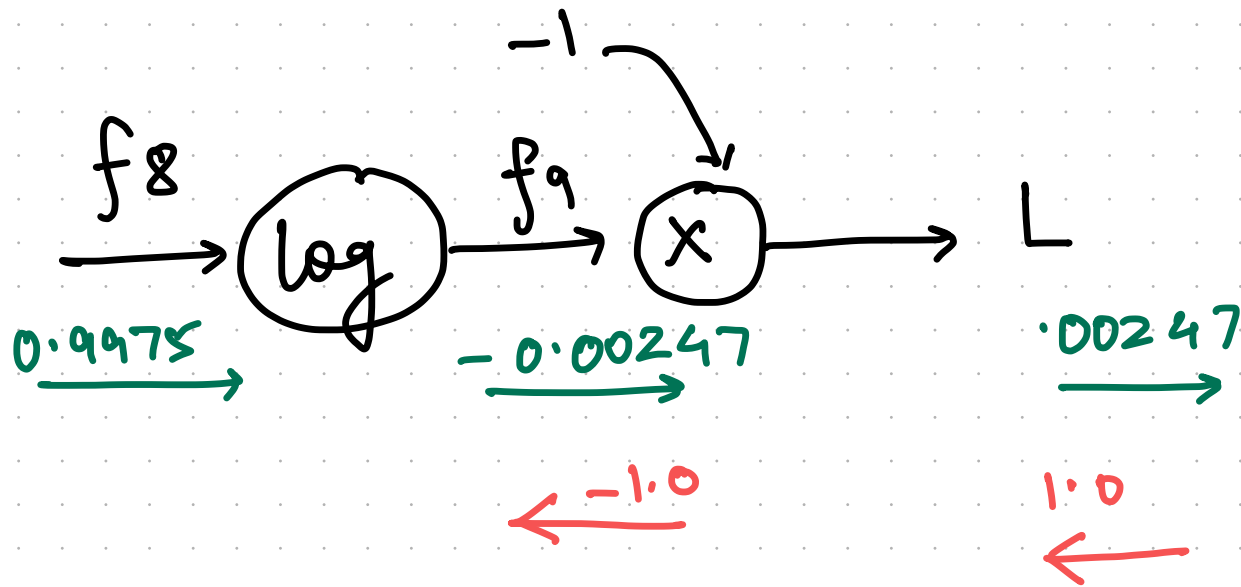
upstream gradient = 1.0

$$L = fg * -1$$

$$\frac{\partial L}{\partial fg} = -1 \qquad \text{LOCAL GRADIENT} = -1$$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$-1$$

$$f_8 \longrightarrow \boxed{\log} \xrightarrow{f_9} \boxed{X} \longrightarrow L$$

$$-0.00247 \longrightarrow$$

$$.00247 \longrightarrow$$

$$\longleftarrow -1.0$$

$$1.0 \longleftarrow$$

$$LOSS = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$-1$$

$$\xrightarrow{f8} \boxed{\log} \xrightarrow{f9} \boxed{X} \longrightarrow L$$

0.9975 →

−0.00247 →

.00247 →

← −1.0

1.0 ←
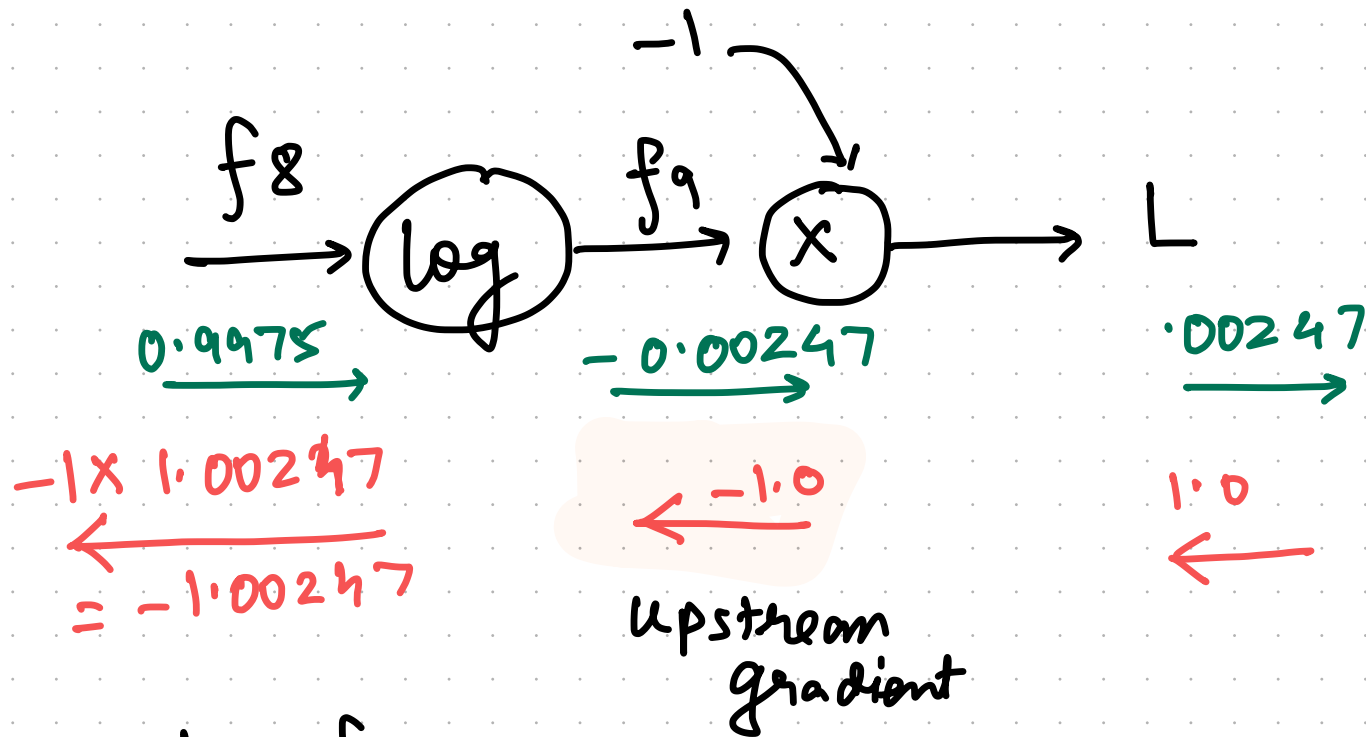
$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$



$$f_9 = \log f_8$$

$$\frac{\partial f_9}{\partial f_8} = \frac{1}{f_8} = \frac{1}{.9975} = 1.00247 = \text{Local gradient}$$

$$LOSS = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$\xrightarrow{f_8} \boxed{\log} \xrightarrow{f_9} \boxed{X} \xrightarrow{} L$$

$-1$

$0.9975$

$-0.00247$

$.00247$

$-1 \times 1.00247$

$-1.0$

$1.0$

$= -1.00247$

upstream gradient

$f_9 = \log f_8$

$$\frac{\partial f_9}{\partial f_8} = \frac{1}{f_8} = \frac{1}{.9975} = 1.00247 = \text{Local gradient}$$

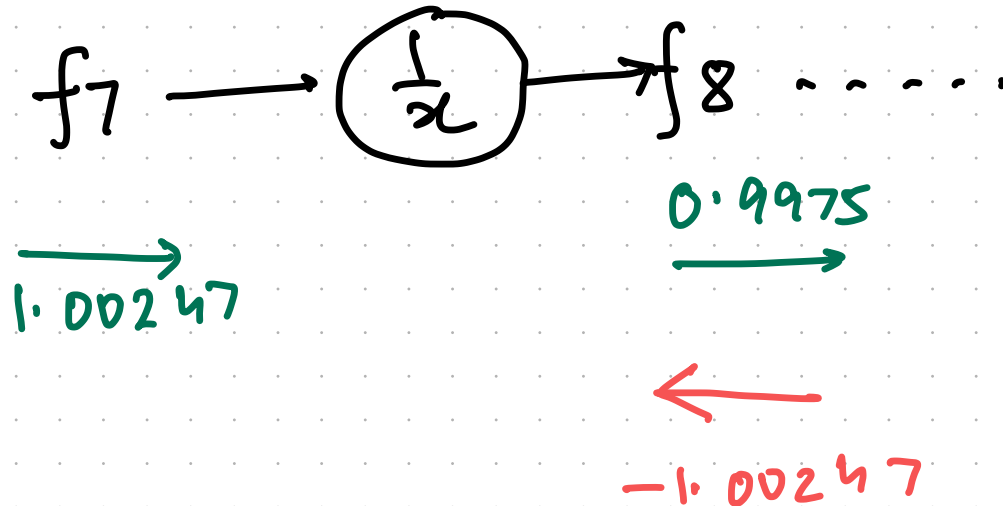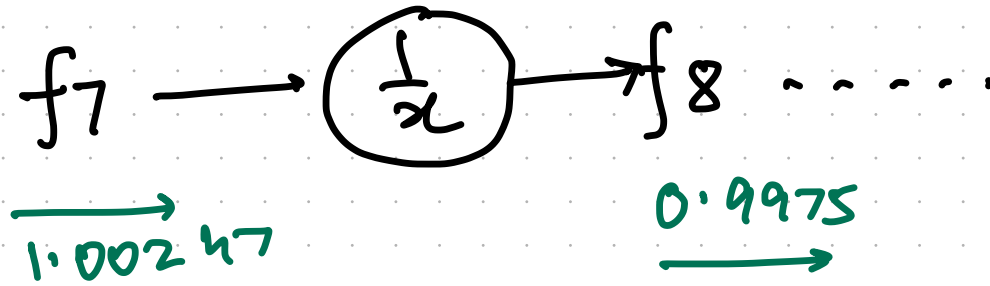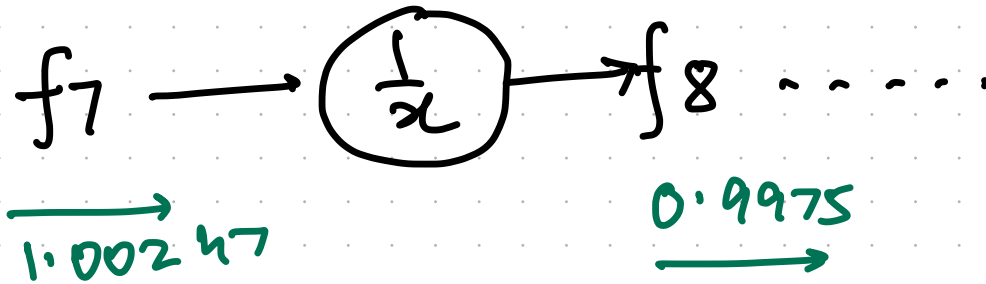$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$f7 \longrightarrow \left(\frac{1}{x}\right) \longrightarrow f8 \cdots \cdots$$

1.00247

0.9975

−1.00247

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$$f_7 \longrightarrow \boxed{\frac{1}{x}} \longrightarrow f_8 \cdot \cdot \cdot \cdot \cdot$$

$\xrightarrow{\text{1.00247}}$

$\xrightarrow{\text{0.9975}}$

$\xleftarrow{\hspace{2cm}}$
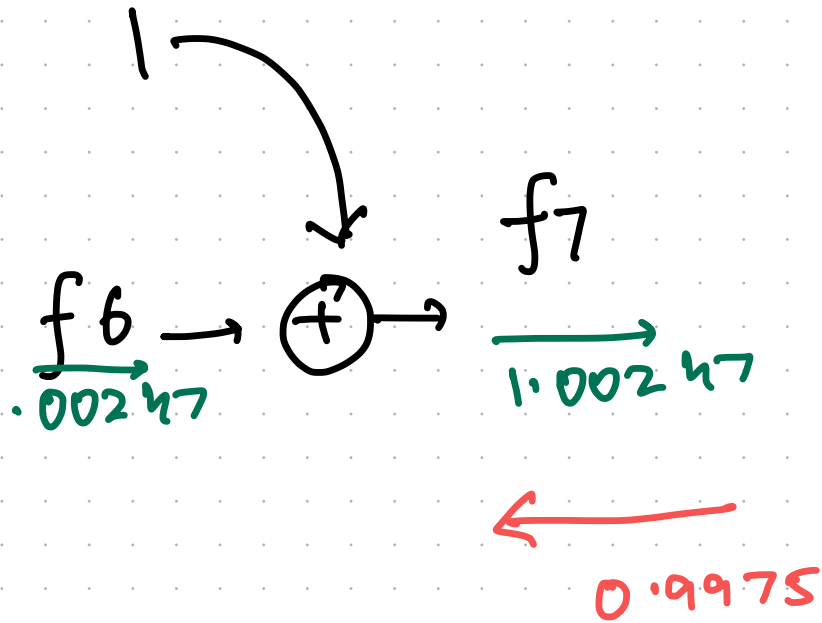
$-1.00247$    upstream gradient

$$f_8 = \frac{1}{f_7} \qquad \frac{\partial f_8}{\partial f_7} = \frac{-1}{f_7^2} = -0.9951$$

$= $ Local gradient

$$LOSS = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$$f_7 \longrightarrow \boxed{\frac{1}{x}} \longrightarrow f_8 \cdots \cdots$$

$$1.00247$$

$$0.9975$$

$$-0.9951 * -1.00247$$
$$= 0.9975$$

$$-1.00247$$

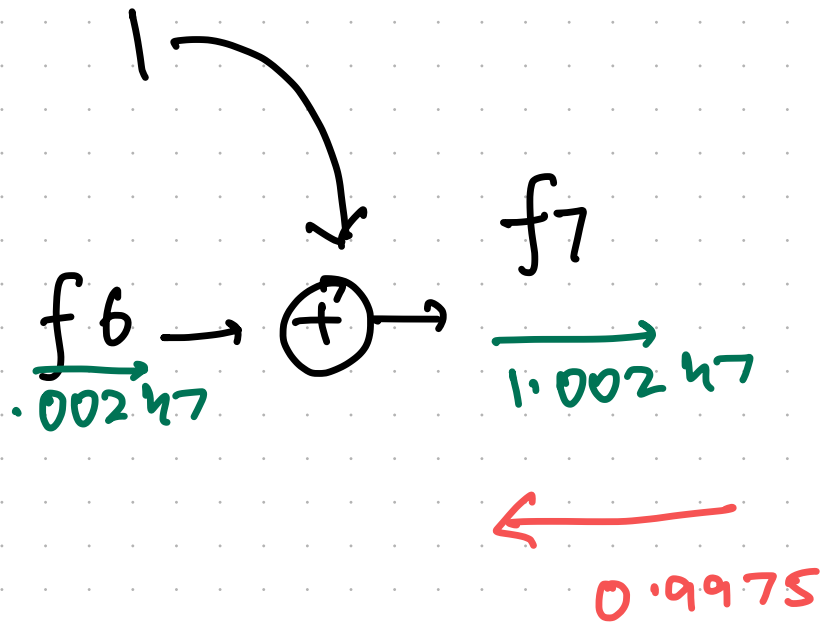upstream gradient

$$f_8 = \frac{1}{f_7} \qquad \frac{\partial f_8}{\partial f_7} = \frac{-1}{f_7^2} = -0.9951$$

= Local gradient

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$1$

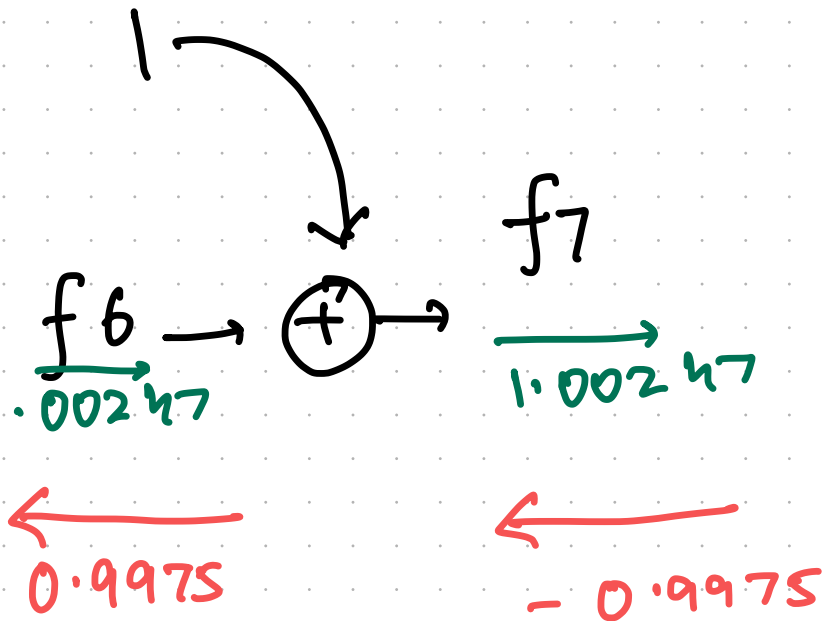$f6 \rightarrow$ $\oplus$ $\rightarrow$ $f7$

$.00247$

$1.00247$

$0.9975$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$1$$

$$f6 \rightarrow \boxed{+} \rightarrow f7$$

$$f6 \xrightarrow{.00247}$$

$$\xrightarrow{1.00247}$$

$$\xleftarrow{0.9975}$$

upstream grad. = 0.9975

local grad. = 1

$$\text{Loss} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

f7

f6 →  (+) →

.00247

1.00247

0.9975

− 0.9975

upstream grad. = 0.9975

local grad. = 1

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

f6

f5 $\longrightarrow$ (exp) $\longrightarrow$ .00247

$-6 \longrightarrow$

$\longleftarrow$

0.9975

upstream grad. $= 0.9975$

local grad. $= \dfrac{\partial f6}{\partial f5} = \dfrac{\partial}{\partial f5} e^{f5} = e^{f5} = e^{-6} = 0.0025$

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$
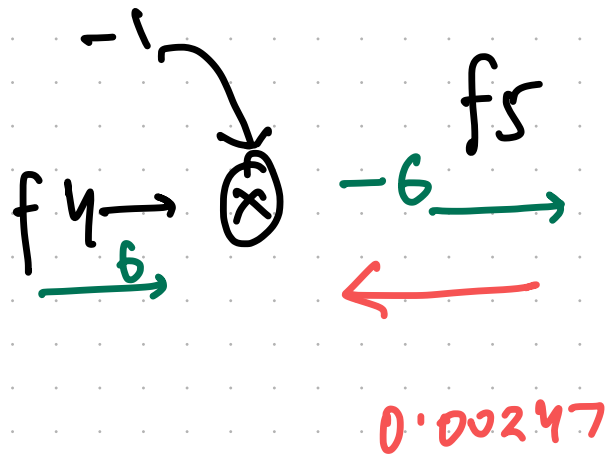
$f5 \longrightarrow$ (exp) $\longrightarrow$ $\overset{f6}{\underset{.00247}{\longrightarrow}}$

$-6 \xrightarrow{\hspace{2cm}}$

$\xleftarrow{\hspace{2cm}} 0.9975$

$\xleftarrow{\hspace{2cm}}$
.0025 * .9975
= 0.00247

upstream grad. = $0.9975$

local grad. $= \dfrac{\partial f6}{\partial f5} = \dfrac{\partial}{\partial f5} e^{f5} = e^{f5} = e^{-6} = 0.0025$

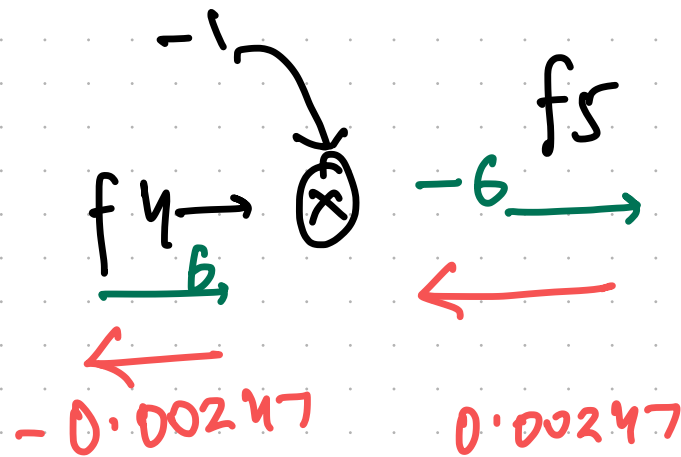$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$$-1 \curvearrowright$$

$$f \, y \longrightarrow \otimes \quad -6 \xrightarrow{f_S}$$

$$f_5$$

$$\xleftarrow{\hspace{2cm}}$$

$$0.00247$$

upstream grad. $= 0.00247$

local grad. $= -1$

$$\text{Loss} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$
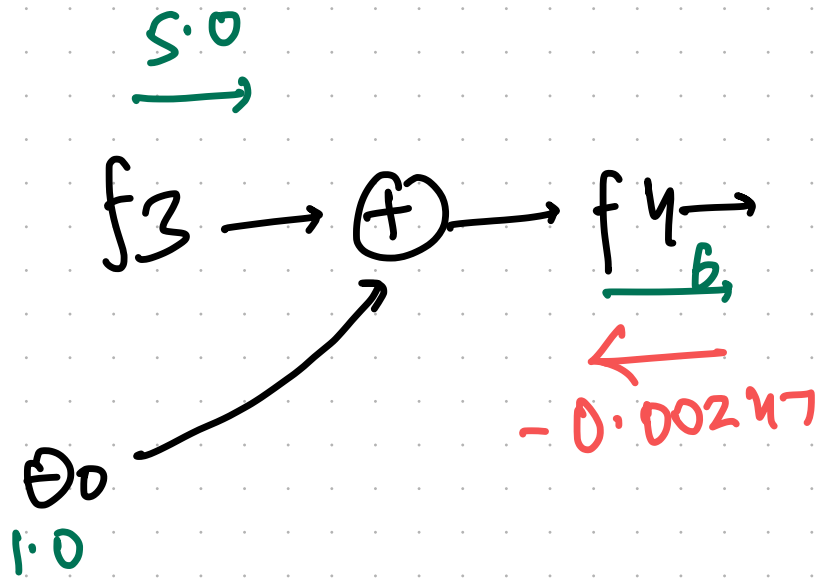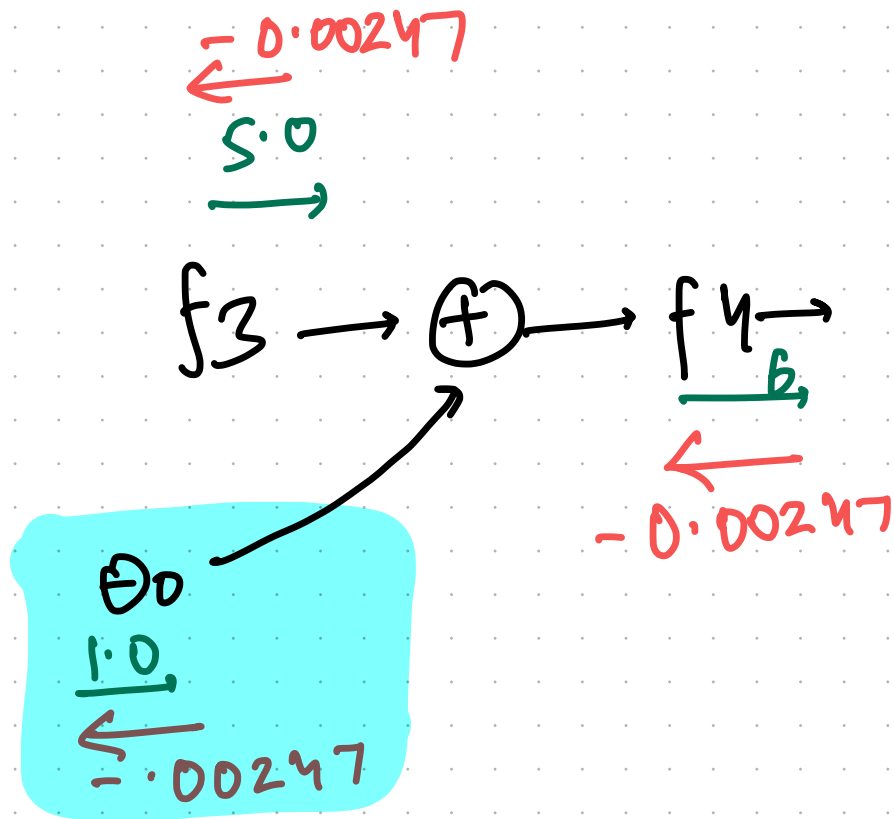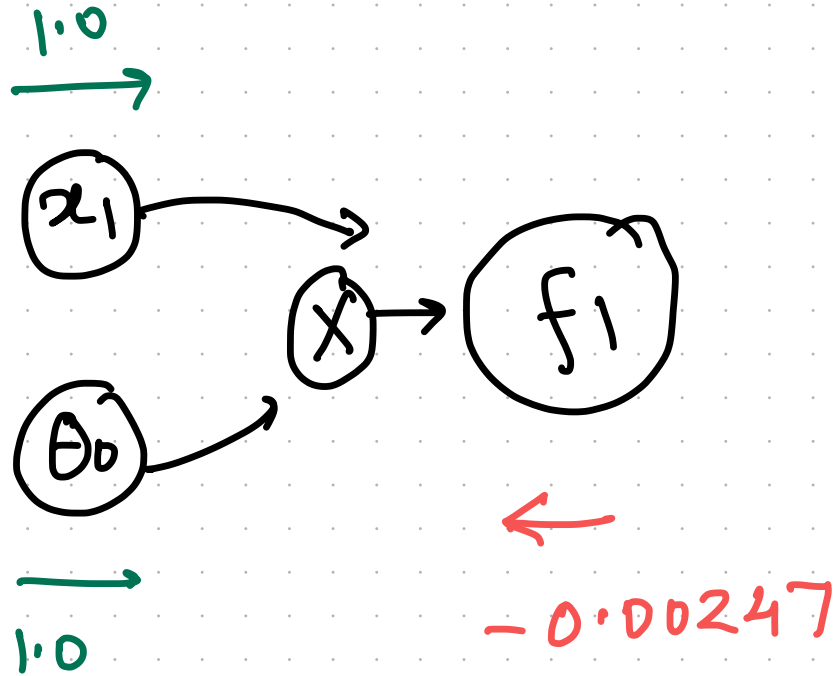


upstream grad. = 0.00247

local grad. = -1

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$



$5.0$

$f3 \longrightarrow (+) \longrightarrow f4 \longrightarrow$
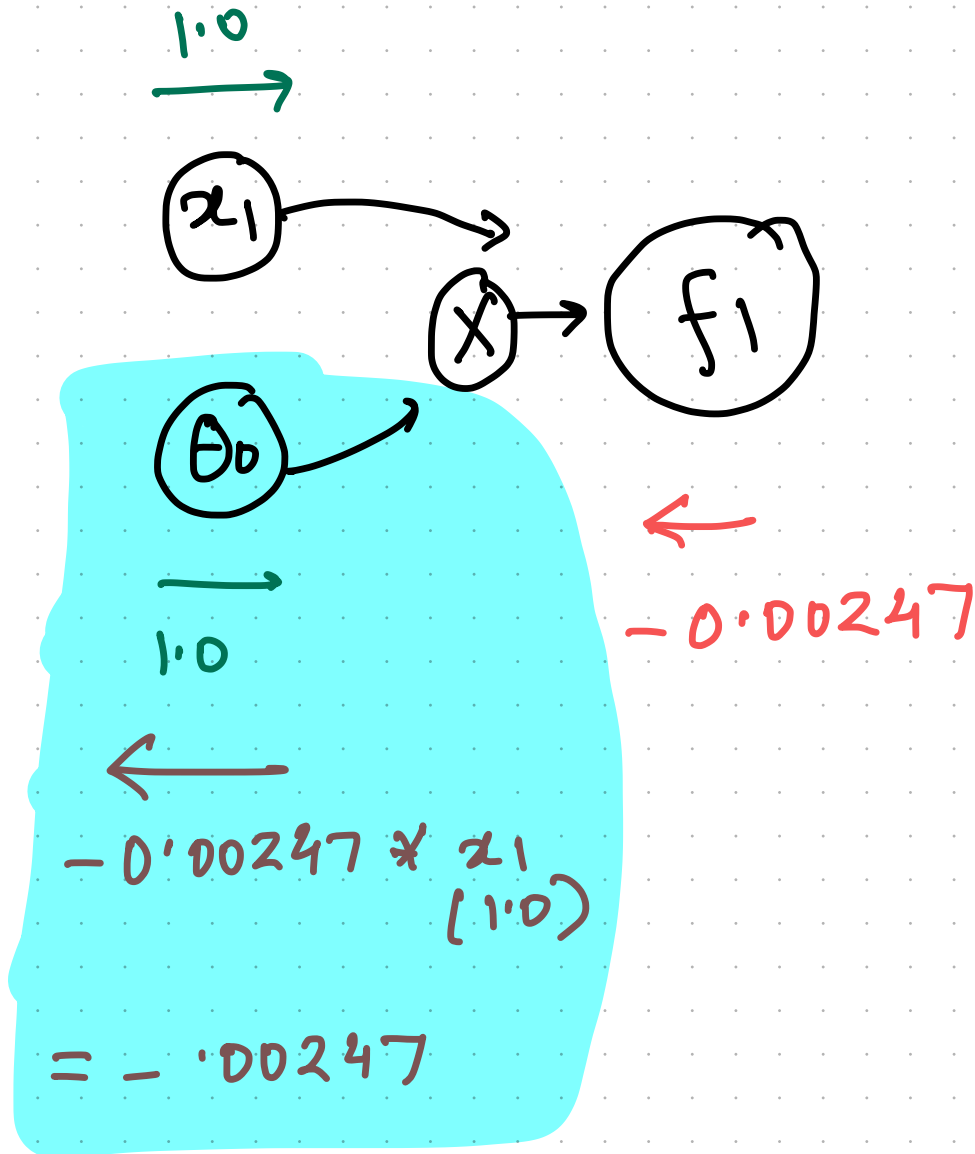
$6,$

$-0.00247$

$\theta_0$

$1.0$

upstream grad. $= -0.00247$

local grad. $(\theta_0) = \dfrac{\partial f4}{\partial \theta_0} = 1$  ;  local grad for $f3 = 1$

$$\text{Loss} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

$= 0.00247$

$5.0$

$f_3 \rightarrow \bigoplus \rightarrow f_4 \rightarrow$

$6,$

$-0.00247$

$\theta_0$

$1.0$

$= .00247$

upstream grad. $= -0.00247$

local grad. $(\theta_0) = \dfrac{\partial f_4}{\partial \theta_0} = 1$ ; local grad for $f_3 = 1$

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$
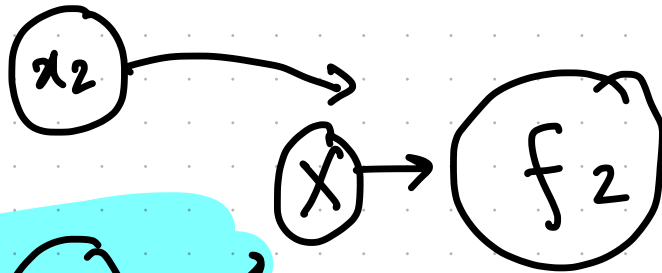
1.0

$x_1$

$\theta_0$

$\times$

$f1$

1.0

−0.00247

$$\text{LOSS} = -1 * \log\left( \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}} \right)$$

1.0 →

$x_1$ → $X$ → $f1$

$\theta_0$

← −0.00247

1.0 →

←

−0.00247 * $x_1$
  (1.0)

= − .00247

$$\text{LOSS} = -1 * \log\left(\frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}}\right)$$

$2.0$

$x_2$

$\times$ → $f_2$

$\theta_2$

$2.0$

$-0.00247$

$= -.00247 * x_2$

$(2.0)$

$= -.0049$

What autodiff library needs to know

(i) $f = a * b$ ; $\dfrac{\partial f}{\partial a} = b$ ; $\dfrac{\partial f}{\partial b} = a$

(ii) $f = a + b$ ; $\dfrac{\partial f}{\partial a} = \dfrac{\partial f}{\partial b} = 1$
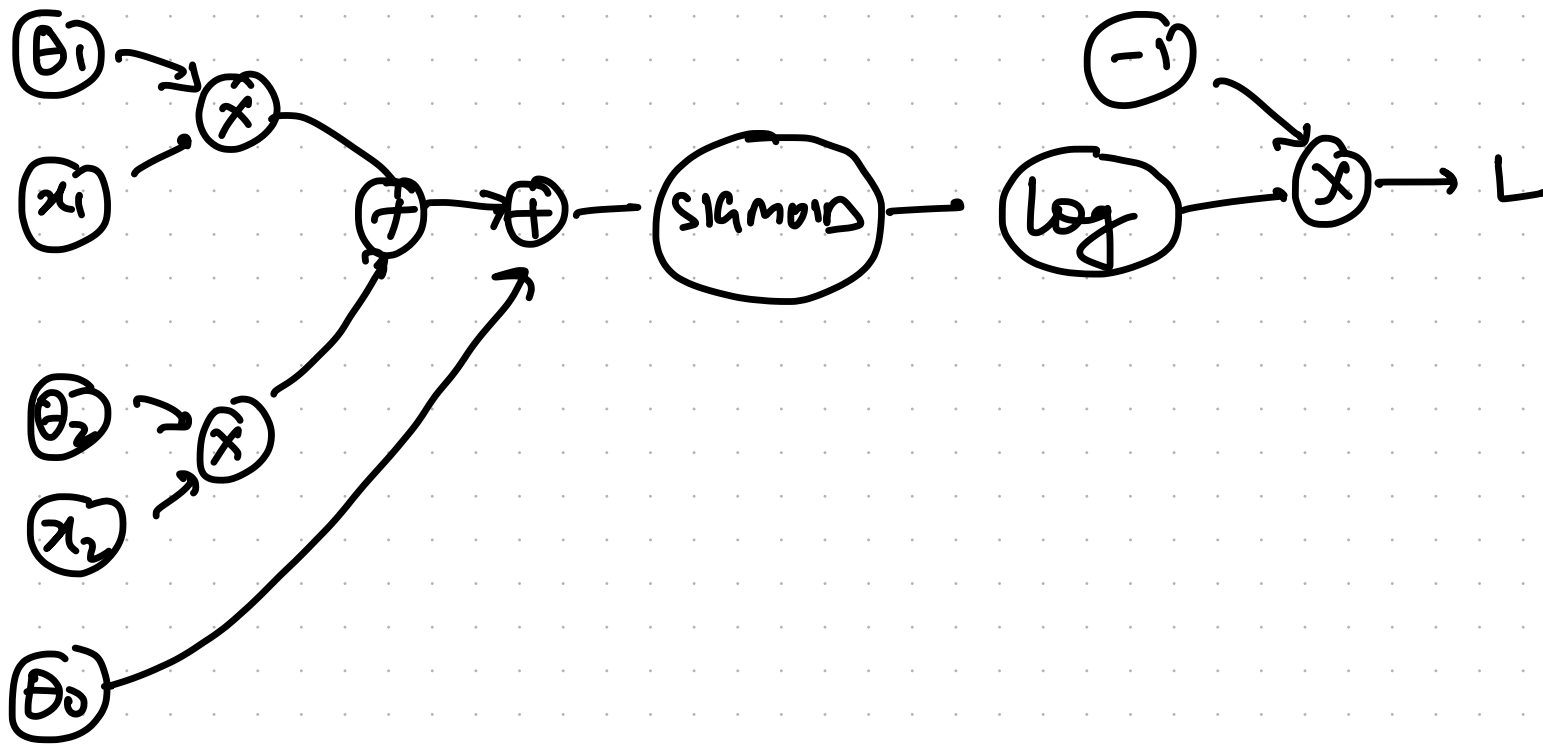
(iii) $f = e^a$ ; $\dfrac{\partial f}{\partial a} = e^a$

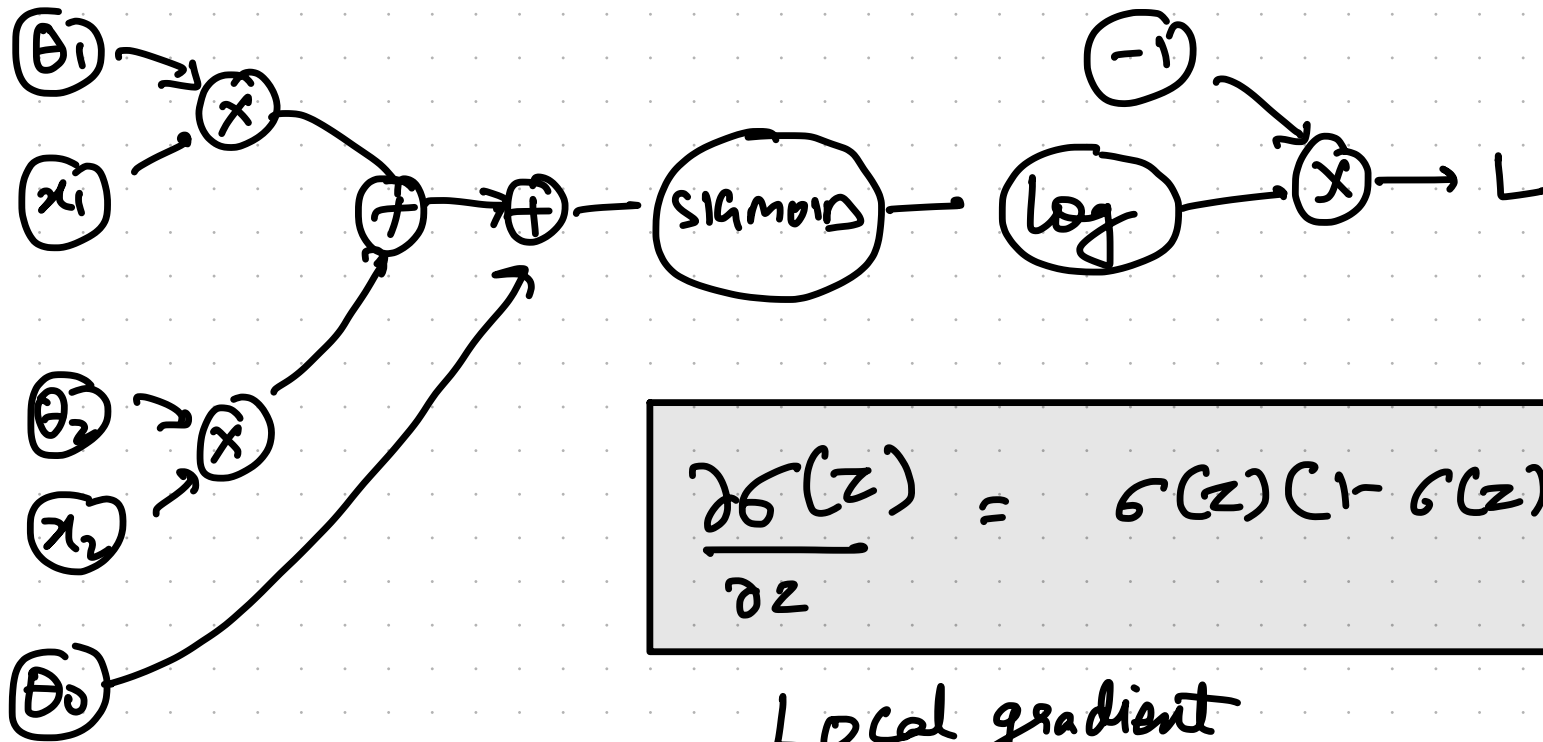(iv) $f = \dfrac{1}{a}$ ; $\dfrac{\partial f}{\partial a} = -1/a^2$

$\vdots$

# Simplifying computation graph

$$L = -1 * \log(\text{SIGMOID}(\theta_0 + \theta_1 x_1 + \theta_2 x))$$

# * Simplifying computation graph

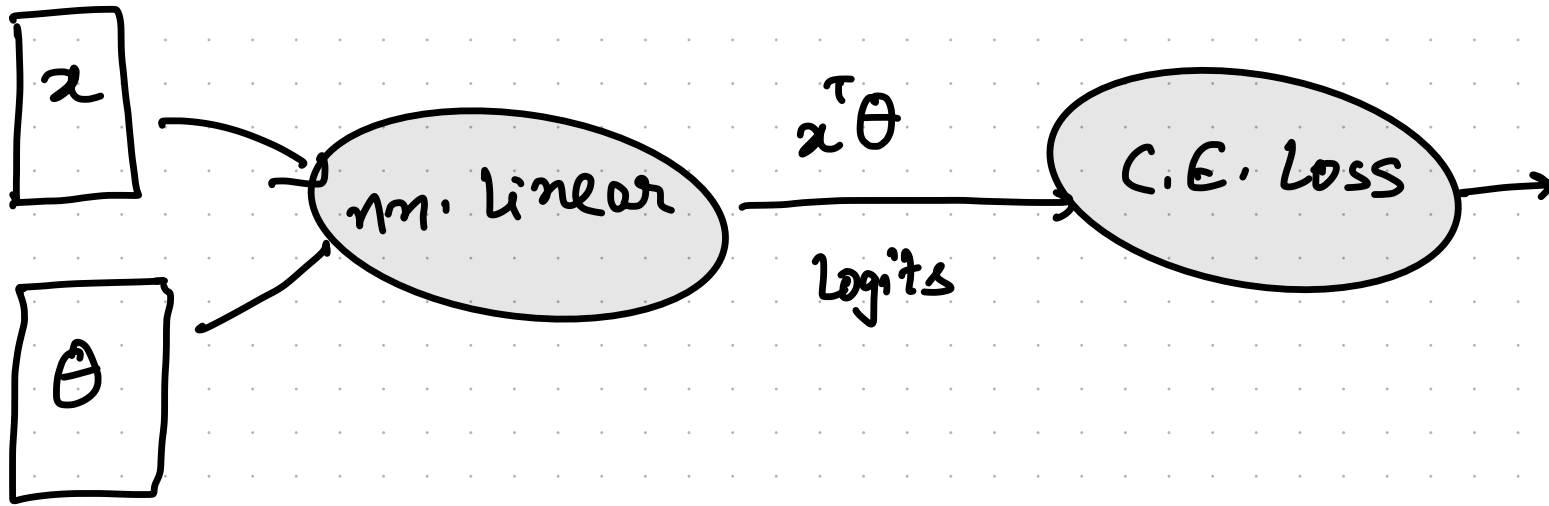$$L = -1 * \log(\text{SIGMOID}(\theta_0 + \theta_1 x_1 + \theta_2 x))$$



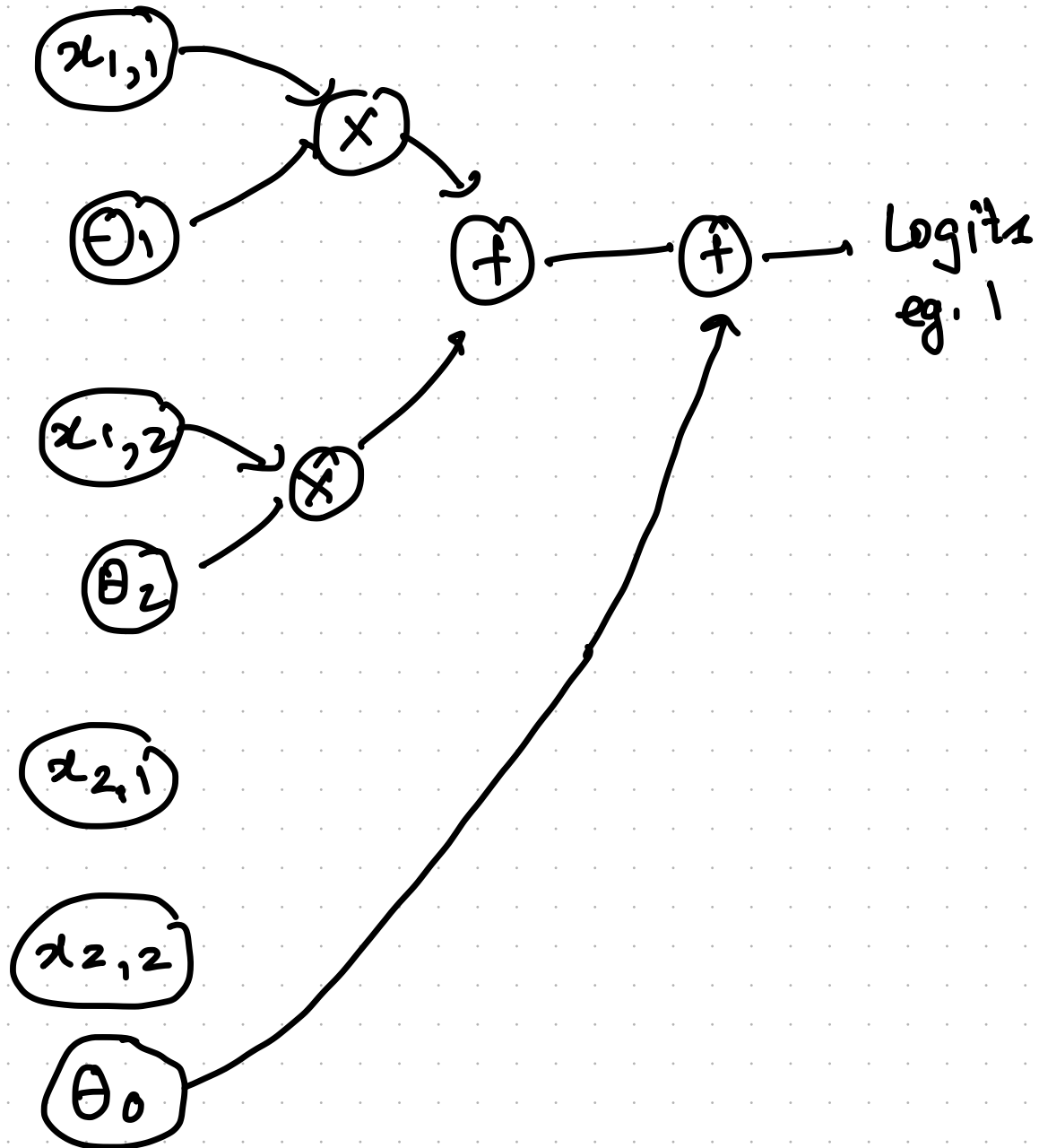$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$

Local gradient

Exercise: show you get same answer
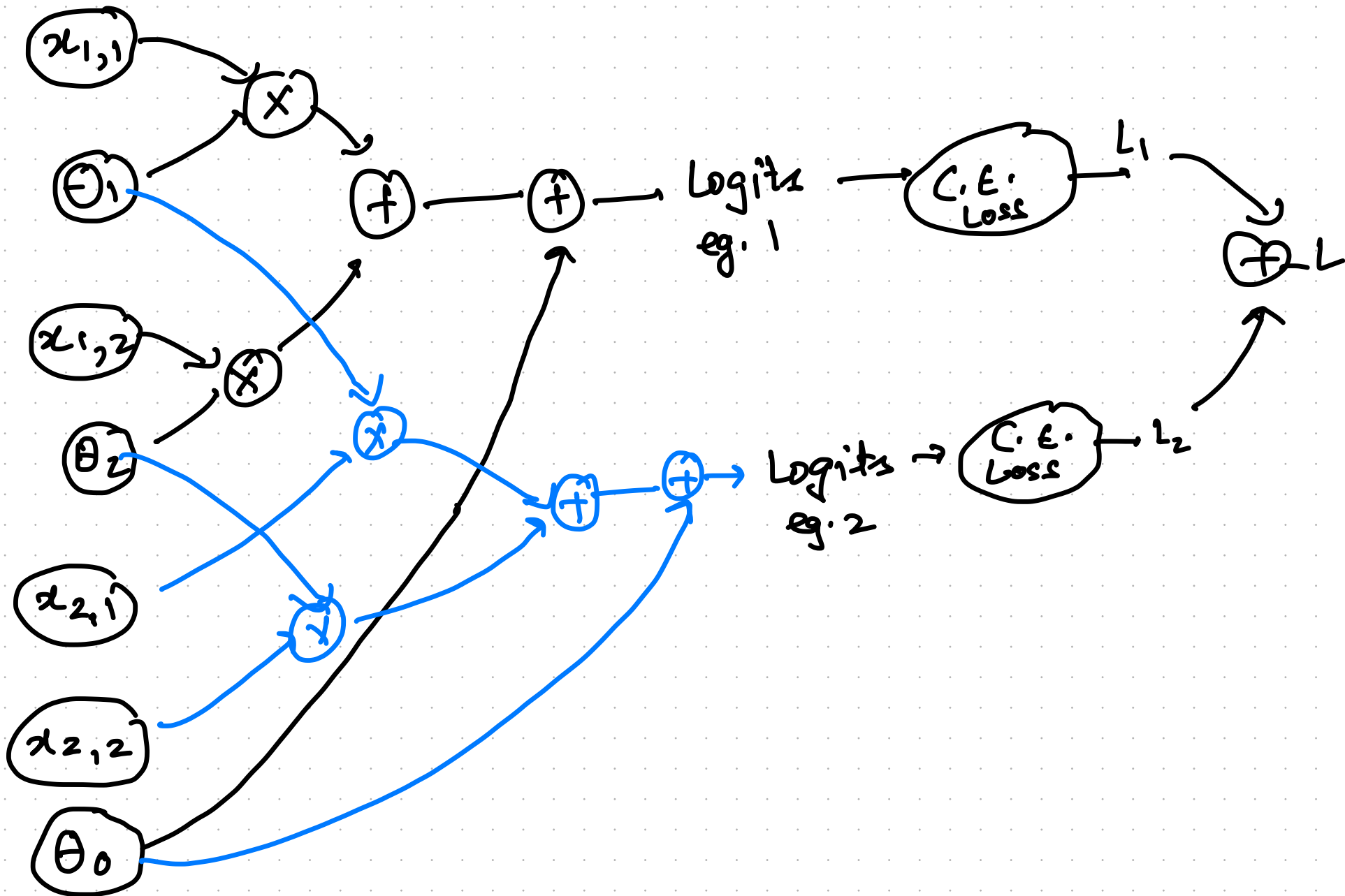as before

**\* Simplifying computation graph**

$$L = -1 * \log(\text{SIGMOID}(\theta_0 + \theta_1 x_1 + \theta_2 x))$$

* Training Over N examples

$x_{1,1}$

$\times$

$\theta_1$

$f$ ——— $f$ ——— Logits
eg. 1

$x_{1,2}$

$\times$

$\theta_2$

$x_{2,1}$

$x_{2,2}$

$\theta_0$

* Training Over N examples

$x_{1,1}$

$\times$

$\theta_1$

$x_{1,2}$

$\times$

$\theta_2$

$+$

$+$

Logits
eg. 1

C.E. Loss

$L_1$

$+ L$

$x_{2,1}$

$x_{2,2}$

$\theta_0$

$\times$

$+$

$+$

Logits →
eg. 2

C.E. Loss

$L_2$

\* Training Over N examples

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

$$\frac{dz}{dt} = \frac{\partial z}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial z}{\partial y} \cdot \frac{dy}{dt},$$

where the ordinary derivatives are evaluated at $t$ and the partial derivatives are evaluated at $(x, y)$.

**\* Training Over N examples**

> 📌 **Chain Rule for One Independent Variable**
>
> Suppose that $x = g(t)$ and $y = h(t)$ are differentiable functions of $t$ and $z = f(x, y)$ is a differentiable function of $x$ and $y$. Then $z = f(x(t), y(t))$ is a differentiable function of $t$ and
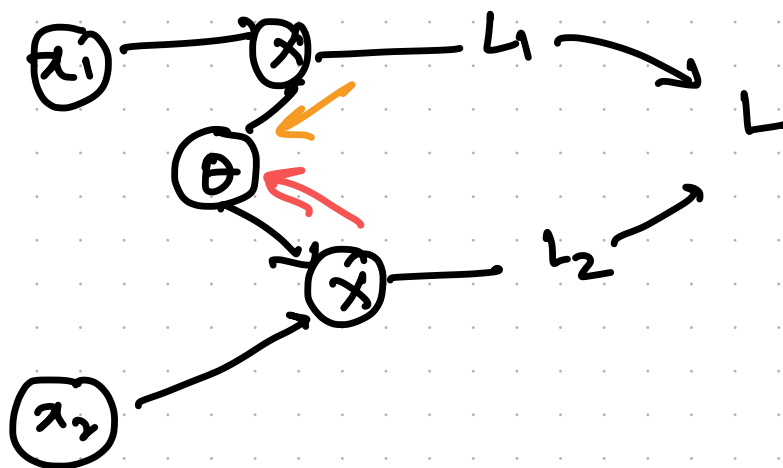>
> $$\frac{dz}{dt} = \frac{\partial z}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial z}{\partial y} \cdot \frac{dy}{dt},$$
>
> where the ordinary derivatives are evaluated at $t$ and the partial derivatives are evaluated at $(x, y)$.

$$L = L_1 + L_2$$

$$L_1 = x_1 \theta$$

$$L_2 = x_2 \theta$$

**\* Training Over N examples**

$$L = L_1 + L_2$$

$$L_1 = x_1 \theta$$

$$L_2 = x_2 \theta$$



$$\frac{\partial L}{\partial \theta} = \longrightarrow + \longrightarrow$$

Addition of all incoming gradients