# Deep Learning

5 Backpropagation-1

Dr. Konda Reddy Mopuri
Dept. of Artificial Intelligence
IIT Hyderabad
Jan-May 2023

# Recap

- Gradient of a scalar valued function $f(\mathbf{x})$: $\mathbf{x} \rightarrow \left( \frac{\partial f}{\partial x_1}, \ldots, \frac{\partial f}{\partial x_D} \right)^T$

# Recap

- Gradient of a scalar valued function $f(\mathbf{x})$: $\mathbf{x} \rightarrow \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_D} \right)^T$

- Gradient of a vector valued function $\mathbf{f}(\mathbf{x})$ is called Jacobian:

$$
\mathbf{J} = \left[ \begin{array}{ccc} \dfrac{\partial \mathbf{f}}{\partial x_1} & \cdots & \dfrac{\partial \mathbf{f}}{\partial x_n} \end{array} \right] = \left[ \begin{array}{c} \nabla^{\mathrm{T}} f_1 \\ \vdots \\ \nabla^{\mathrm{T}} f_m \end{array} \right] = \left[ \begin{array}{ccc} \dfrac{\partial f_1}{\partial x_1} & \cdots & \dfrac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial f_m}{\partial x_1} & \cdots & \dfrac{\partial f_m}{\partial x_n} \end{array} \right]
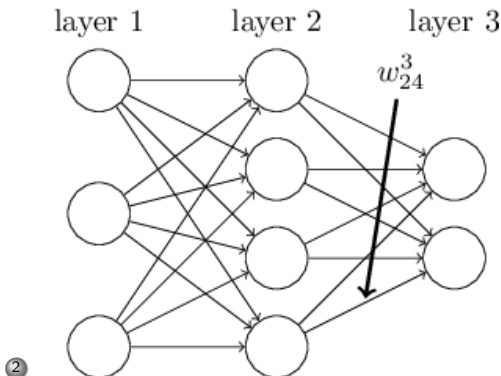$$

# MLP: Some Notation

1. $w_{jk}^{l}$ is the weight connecting $j^{th}$ neuron in $l^{th}$ layer and $k^{th}$ neuron in $(l-1)^{st}$ layer

# MLP: Some Notation

1. $w_{jk}^l$ is the weight connecting $j^{th}$ neuron in $l^{th}$ layer and $k^{th}$ neuron in $(l-1)^{st}$ layer

# MLP: Some Notation

1. $b_j^l$ is the bias of $j^{th}$ neuron in $l^{th}$ layer

# MLP: Some Notation

1. $b_j^l$ is the bias of $j^{th}$ neuron in $l^{th}$ layer
2. $x_j^l$ is the activation (output) of $j^{th}$ neuron in $l^{th}$ layer

# MLP: Some Notation

1. $b_j^l$ is the bias of $j^{th}$ neuron in $l^{th}$ layer
2. $x_j^l$ is the activation (output) of $j^{th}$ neuron in $l^{th}$ layer
3.

$$x_j^l = \sigma\left(\sum_k w_{jk}^l x_k^{l-1} + b_j^l\right)$$

# MLP: Some Notation

1. $b_j^l$ is the bias of $j^{th}$ neuron in $l^{th}$ layer
2. $x_j^l$ is the activation (output) of $j^{th}$ neuron in $l^{th}$ layer
3.

$$x_j^l = \sigma\left(\sum_k w_{jk}^l x_k^{l-1} + b_j^l\right)$$

4. Vector of activations (or, biases) at a layer $l$ is denoted by a bold-faced $\mathbf{x}^l$ ( or $\mathbf{b}^l$) and $W^l$ is the matrix of weights into layer $l$

# MLP: Some Notation

1. $s_j^l$ is the weighted input to $j^{th}$ neuron in $l^{th}$ layer

# MLP: Some Notation

1. $s_j^l$ is the weighted input to $j^{th}$ neuron in $l^{th}$ layer
2. $s_j^l = \sum_k w_{jk}^l x_k^{l-1} + b_j^l$

# MLP: Some Notation

1.  $s_j^l$ is the weighted input to $j^{th}$ neuron in $l^{th}$ layer
2.  $s_j^l = \sum_k w_{jk}^l x_k^{l-1} + b_j^l$
3.  $\mathbf{s}^l = W^l \mathbf{x}^{l-1} + \mathbf{b}^l$

# MLP: Some Notation

1. $s_j^l$ is the weighted input to $j^{th}$ neuron in $l^{th}$ layer
2. $s_j^l = \sum_k w_{jk}^l x_k^{l-1} + b_j^l$
3. $\mathbf{s}^l = W^l \mathbf{x}^{l-1} + \mathbf{b}^l$
4. $\sigma$ is the activation function that applies element-wise

# Gradient descent on MLP

- Loss is $\mathcal{L}(W, \mathbf{b}) = \sum_n l(f(x_n; W, \mathbf{b}), y_n) = \sum_n l(\mathbf{x}^L, y_n)$ ($L$ is the number of layers in the MLP)

# Gradient descent on MLP

- Loss is $\mathcal{L}(W, \mathbf{b}) = \sum_n l(f(x_n; W, \mathbf{b}), y_n) = \sum_n l(\mathbf{x}^L, y_n)$ ($L$ is the number of layers in the MLP)

- For applying Gradient descent, we need gradient of individual sample loss with respect to all the model parameters

$$l_n = l(f(x_n; W, \mathbf{b}), y_n)$$

$\frac{\partial l_n}{\partial W_{jk}^{(l)}}$ and $\frac{\partial l_n}{\partial \mathbf{b}_j^{(l)}}$ for all layers $l$

## Forward pass operation

$$x^{(0)} = x \xrightarrow{W^{(1)},\mathbf{b}^{(1)}} s^{(1)} \xrightarrow{\sigma} x^{(1)} \xrightarrow{W^{(2)},\mathbf{b}^{(2)}} s^{(2)} \ldots x^{(L-1)} \xrightarrow{W^{(L)},\mathbf{b}^{(L)}} s^{(L)} \xrightarrow{\sigma} x^{(L)} = f(x; W, \mathbf{b})$$

Formally, $x^{(0)} = x$, $f(x; W, \mathbf{b}) = x^{(L)}$

$$\forall l = 1, \ldots, L \quad \begin{cases} s^{(l)} &= W^{(l)} x^{(l-1)} + \mathbf{b}^{(l)} \\ x^{(l)} &= \sigma(s^{(l)}) \end{cases}$$

# Chain rule of differential calculus

- Core concept of backpropagation

# Chain rule of differential calculus

- Core concept of backpropagation
-
$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

# Chain rule of differential calculus

- Core concept of backpropagation

-
$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x)$$

-
$$\frac{\partial}{\partial x} g(f(x)) = \left. \frac{\partial g(a)}{\partial a} \right|_{a=f(x)} \cdot \frac{\partial f(x)}{\partial x}$$
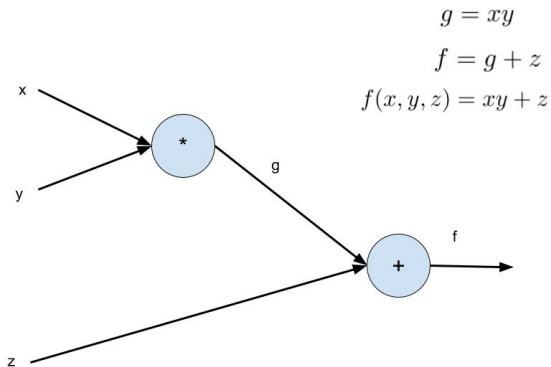
# Chain rule of differential calculus



The Chain Rule

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

$$\frac{dy}{dx} = \left( \begin{array}{c} \text{Differentiate} \\ \text{outer function} \\ \\ \text{Keep the inside} \\ \text{the same} \end{array} \right) \left( \begin{array}{c} \text{Differentiate} \\ \text{inner function} \end{array} \right)$$
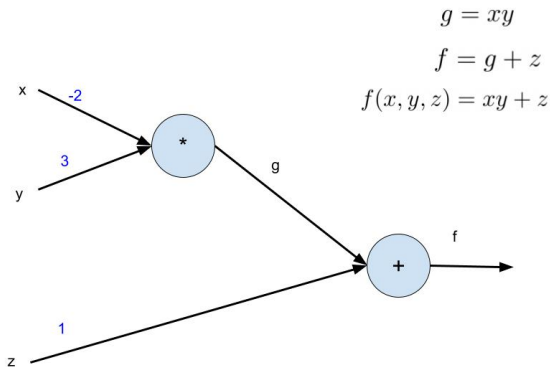
# Chain rule of differential calculus

1. $f(x) = e^{sin(x^2)}$, let's find $\frac{\partial f}{\partial x}$ (work it out on the board)

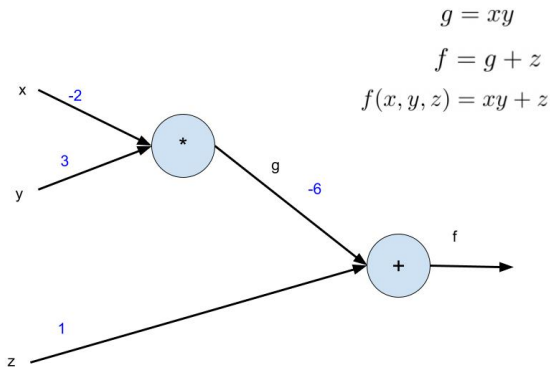# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$L = (f - t)^2$$
$$t = -4 \quad \text{(for this input)}$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$L = (f - t)^2$$
$$t = -4 \quad \text{(for this input)}$$
$$\frac{\partial L}{\partial f} = 2(f - t)$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial f} = -2$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

x    -2

3

y

\*

g    -6

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

+

f    -5

$$\frac{\partial L}{\partial f} = -2$$

1

z

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

x  -2

3

y

$$\frac{\partial L}{\partial g} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial g}$$

g  -6

f  -5

$$\frac{\partial L}{\partial f} = -2$$

1

z

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

x  -2

3

y

$*$

g  -6

$$\frac{\partial L}{\partial g} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial g}$$

$+$

f  -5

$$\frac{\partial L}{\partial f} = -2$$

1

z

# Chain rule of differential calculus



$$\frac{\partial f}{\partial g} = 1 \longleftarrow \begin{array}{l} g = xy \\ f = g + z \\ f(x, y, z) = xy + z \end{array}$$

$$\frac{\partial L}{\partial g} = \frac{\partial L}{\partial f}\frac{\partial f}{\partial g}$$

$$\frac{\partial L}{\partial f} = -2$$

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

x  -2

3

y

*

g  -6

$$\frac{\partial L}{\partial g} = -2$$

f  -5

+

$$\frac{\partial L}{\partial f} = -2$$

1

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

z

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial g} = -2$$

$$\frac{\partial L}{\partial f} = -2$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial z}$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial g} = -2$$

$$\frac{\partial L}{\partial f} = -2$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial f}\frac{\partial f}{\partial z}$$

# Chain rule of differential calculus



$$\frac{\partial f}{\partial z} = 1 \longleftarrow \begin{array}{l} g = xy \\ f = g + z \end{array}$$

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

$$\frac{\partial L}{\partial g} = -2$$

$$\frac{\partial L}{\partial f} = -2$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial z}$$

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

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

x   -2

3

y

*

g   -6

$\frac{\partial L}{\partial g} = -2$

f   -5

+

$\frac{\partial L}{\partial f} = -2$

1

$\frac{\partial L}{\partial z} = -2$

z

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial g}\frac{\partial g}{\partial x}$$

x   -2

3

y

* (node)

g   -6

$$\boxed{\frac{\partial L}{\partial g} = -2}$$

1

$$\boxed{\frac{\partial L}{\partial z} = -2}$$

z

+ (node)

f   -5

$$\boxed{\frac{\partial L}{\partial f} = -2}$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial g} \frac{\partial g}{\partial x}$$

x    -2

3

y

g    -6

$$\frac{\partial L}{\partial g} = -2$$

*

+

f    -5

$$\frac{\partial L}{\partial f} = -2$$

1

$$\frac{\partial L}{\partial z} = -2$$

z

# Chain rule of differential calculus



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial g}\frac{\partial g}{\partial x}$$

$$\frac{\partial g}{\partial x} = y \quad \longleftarrow \quad g = xy$$

$$f = g + z$$

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

x  -2

3

y

g  -6

$$\frac{\partial L}{\partial g} = -2$$

*

+

f  -5

$$\frac{\partial L}{\partial f} = -2$$

1

$$\frac{\partial L}{\partial z} = -2$$

z

# Chain rule of differential calculus

$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial x} = -6$$

x    -2

3

y

*

g    -6

$$\frac{\partial L}{\partial g} = -2$$

f    -5

+

$$\frac{\partial L}{\partial f} = -2$$

1

$$\frac{\partial L}{\partial z} = -2$$

z

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial x} = -6$$

x   -2

3

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

g   -6

$$\frac{\partial L}{\partial g} = -2$$

f   -5

$$\frac{\partial L}{\partial f} = -2$$

1

z

$$\frac{\partial L}{\partial z} = -2$$

# Chain rule of differential calculus



$$g = xy$$
$$f = g + z$$
$$f(x, y, z) = xy + z$$

$$\frac{\partial L}{\partial x} = -6$$

x $\quad$ -2

y $\quad$ 3

$$\frac{\partial L}{\partial y} = 4$$

*

g $\quad$ -6

$$\frac{\partial L}{\partial g} = -2$$

f $\quad$ -5

+

$$\frac{\partial L}{\partial f} = -2$$

z $\quad$ 1

$$\frac{\partial L}{\partial z} = -2$$

# Gradient Flow

# Gradient Flow



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial g}\frac{\partial g}{\partial x}$$

x

$\frac{\partial L}{\partial g}$

*

g

y

# Gradient Flow



$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial g}\frac{\partial g}{\partial x}$$

x

*

$\frac{\partial L}{\partial g}$

g

y

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial q}\frac{\partial g}{\partial x}$$

x

Local Gradient

Upstream Gradient

$$\frac{\partial L}{\partial q}$$
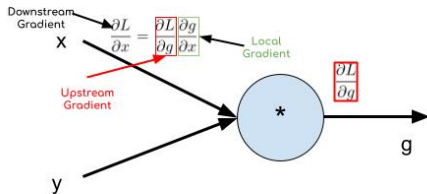
g

y

# Gradient Flow

# Chain rule of differential calculus for an MLP

○

$$J_{f_N \circ f_{N-1} \circ \ldots f_1(x)} = J_{f_N(f_{N-1}(\ldots f_1(x)))} \cdot J_{f_{N-1}(f_{N-2}(\ldots f_1(x)))} \cdot \ldots \cdot J_{f_2(f_1(x))} \cdot J_{f_1(x)}$$

$J_{f(x)}$ is Jacobian of $f$ computed at x.