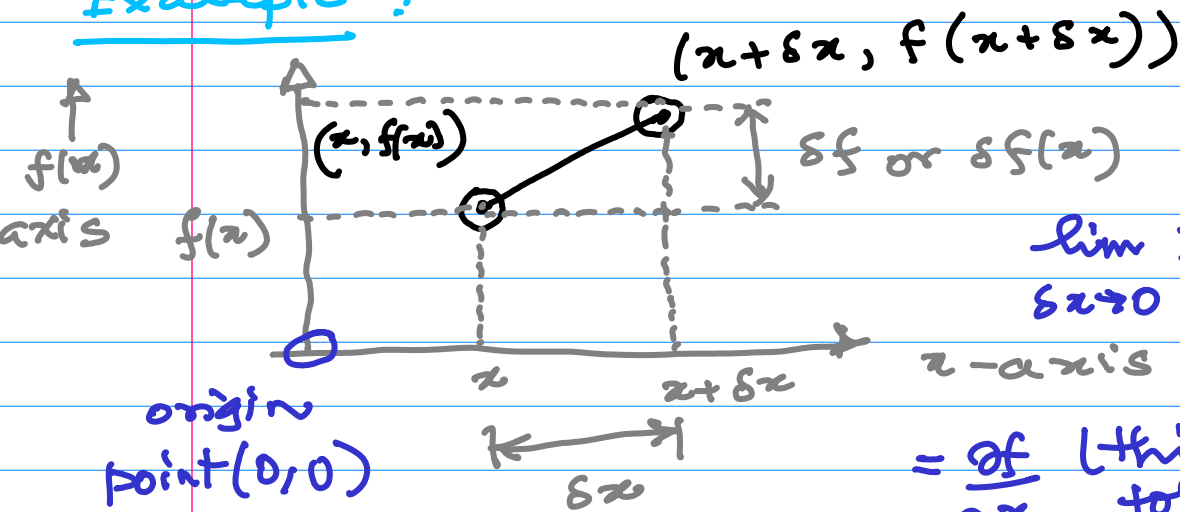


DERIVATIVES

The most general derivative is NOT the total derivative, but the partial derivatives.

[1-D] one independent variable $x \rightarrow$ scalar
 one dependent variable $f(x) \rightarrow$ scalar function

Example : e.g., audio signal $f(t)$



$$\lim_{\delta x \rightarrow 0} \frac{f(x + \delta x) - f(x)}{(x + \delta x) - (x)}$$

$$= \frac{\partial f}{\partial x} \quad (\text{this is also the total derivative } df/dx)$$

Why? 1 scalar variable.

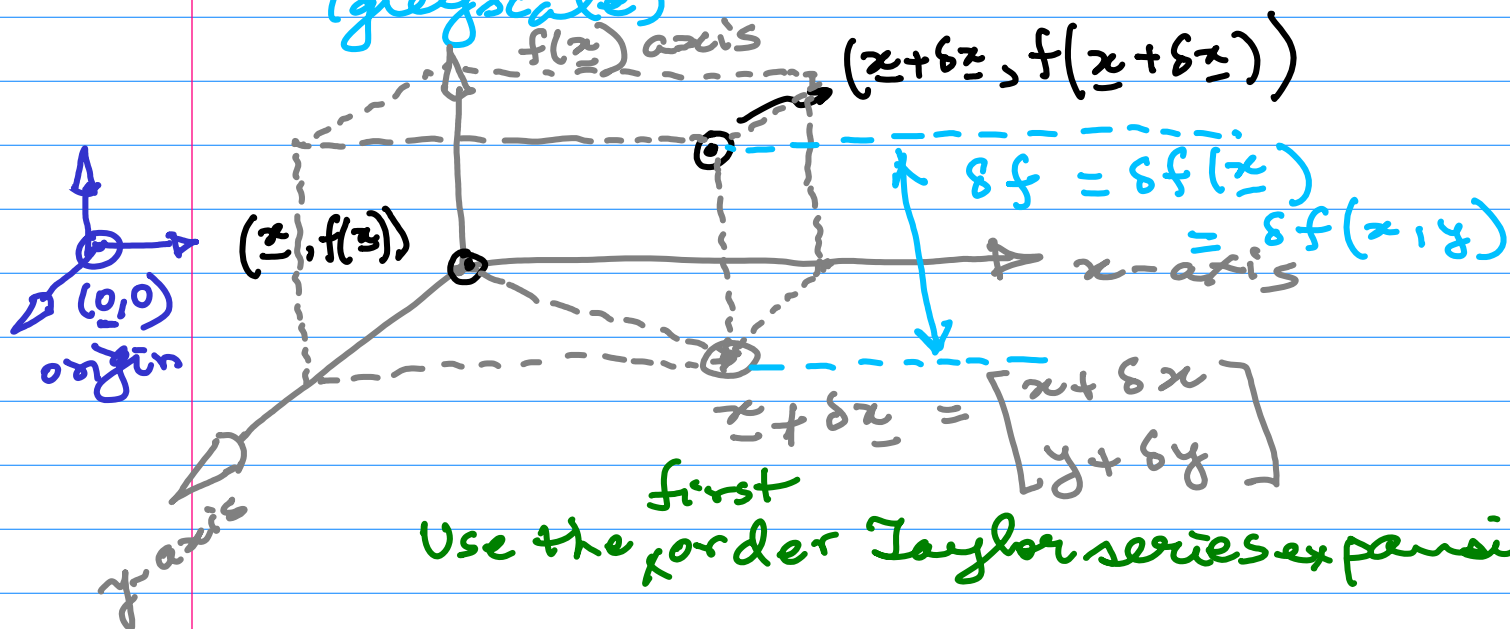
$$\lim_{\delta x \rightarrow 0} : \underbrace{f(x + \delta x) - f(x)}_{\delta f \text{ or } \delta f(x)} = \left(\frac{\partial f}{\partial x} \right) \underbrace{\delta x}_{\text{small change in the independent variable}}$$

Small change in the dependent variable, or the function

the independent variable.

[2-D] two independent variables $\underline{x} = \begin{bmatrix} x \\ y \end{bmatrix} \rightarrow$ vector
 one dependent variable $f(\underline{x}) \rightarrow$ scalar.

Example $I(x, y)$ or $I(\underline{x})$
 the intensity of a pixel at position $\underline{x} = \begin{bmatrix} x \\ y \end{bmatrix}$
 (greyscale)



first
 Use the Taylor series expansion

$$f(\underline{x} + \delta \underline{x}) = f(\underline{x}) + \nabla f \cdot \delta \underline{x}$$

$$\underbrace{f(\underline{x} + \delta \underline{x}) - f(\underline{x})}_{\substack{\delta f \text{ or } \delta f(\underline{x}) \\ \text{or } \delta f(x, y)}} = \underbrace{\nabla f}_{\begin{bmatrix} \partial f / \partial x \\ \partial f / \partial y \end{bmatrix}} \cdot \underbrace{\delta \underline{x}}_{\begin{bmatrix} \delta x \\ \delta y \end{bmatrix}}$$

$$\underbrace{\frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y}_{= \sum_{\text{var}: x, y} \left(\frac{\partial f}{\partial \text{var}} \right) (\delta \text{var})}$$

3-D Three independent variables $\underline{x} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$ ↗ vector
 one dependent variable $f(\underline{x})$ ↗ scalar

Example: video (greyscale) $I(x, y, t)$: Intensity at pixel (x, y) in frame # t .

Impossible to visualise $I(x, y, t) \Rightarrow 4$ -D entity
 Use the first order Taylor series approximation.

$$f(\underline{x} + \delta \underline{x}) = f(\underline{x}) + \bar{\nabla} f \cdot \delta \underline{x}$$

$$\underbrace{f(\underline{x} + \delta \underline{x}) - f(\underline{x})}_{\substack{\delta f \text{ or } \delta f(\underline{x}) \\ \text{or } \delta f(x, y, z)}} = \underbrace{\begin{bmatrix} \partial f / \partial x \\ \partial f / \partial y \\ \partial f / \partial z \end{bmatrix}}_{\sim} \cdot \begin{bmatrix} \delta x \\ \delta y \\ \delta z \end{bmatrix}$$

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

$$= \sum_{\text{var} = x, y, z} \left(\frac{\partial f}{\partial \text{var}} \right) (\delta \text{var})$$

Generalise to a function of D variables

$$f(\underline{x}), \quad \underline{x} = \begin{bmatrix} x_D \\ \vdots \\ x_1 \end{bmatrix} \quad \text{1st order Taylor series expansion}$$

$$\left. \begin{array}{l} \delta f \text{ or } \delta f(\underline{x}) \\ \text{or, } \delta f(x_1 \dots x_D) \end{array} \right\} = \bar{\nabla} f \cdot \delta \underline{x} = \sum_{i=1}^D \frac{\partial f}{\partial x_i} \delta x_i$$

Take-home point: The total change is always a summation

$$\delta f = \delta f(\underline{x}) = \sum_{i=1}^D \frac{\partial f}{\partial x_i} \delta x_i$$

Consider another variable t

(all the x_i 's are functions of this variable t)

$$\frac{\delta f}{\delta t} = \sum_{i=1}^D \frac{\partial f}{\partial x_i} \frac{\delta x_i}{\delta t}$$

We can take the limit as $\delta t \rightarrow 0$

$$\frac{\partial f}{\partial t} = \sum_{i=1}^D \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial t}$$

Here, we had one variable t , so the partial derivative is also the total derivative.

$$\frac{df}{dt} = \sum_{i=1}^D \frac{\partial f}{\partial x_i} \frac{dx_i}{dt}$$

Now, consider a set of variables $t_j : j \in \{1, D\}$

(all the x_i 's are functions of t_j)

$i \in \{1, D\}$

$j \in \{1, D\}$

$$\forall j : \frac{\partial f}{\partial t_j} = \sum_{i=1}^D \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial t_j} \quad \text{Chain Rule}$$

Compact Moral of the story: if f depends on many x_i , then

$$\frac{\partial f}{\partial t} = \sum_{\forall i} \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial t}$$

Perceptron & MLP: some closing notes

- (*) Key difference: MLP and the perceptron:
MLP uses continuous sigmoidal non-linearities in the hidden units, whereas the perceptron uses a step function non-linearities
- (*) Variants: skip layers: either direct connection, or with a small first layer weight (so that over its operating range, the hidden unit is effectively linear), compensating with a large weight value from the hidden unit to the output
- (*) Sparse network (CNN)

Math overview / recap: → 1st order Taylor series approximation

$$\underbrace{E(\underline{w} + \delta \underline{w})}_{\substack{\text{error function} \\ \text{minimise}}} = E(\underline{w}) + \underbrace{(\nabla E)}_{\substack{\text{weights} \\ \text{(All)}}} \cdot (\delta \underline{w})$$

Overall aim to find a weight vector, which minimises an error function $E(\underline{w})$

At the extremum, $\nabla E = \underline{0}$ (scalar) (vector)

max/min/saddle point vector

$$\begin{bmatrix} \partial E / \partial w_2 \\ \partial E / \partial w_1 \end{bmatrix}_{(2-D)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Local Quadratic Approximation

$$E(\underline{w} + \delta \underline{w}) = E(\underline{w}) + \underbrace{\nabla E \cdot (\delta \underline{w})}_{(\delta \underline{w})^T \nabla E} + \frac{1}{2} (\delta \underline{w})^T H (\delta \underline{w}) + \text{higher order terms}$$

↙ Hessian

Example: first order

$$I(x+p, y+q, t+v) = I(x, y, t) + p I_x + q I_y + v I_t$$

Second order term

$$H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j} \Big|_{\underline{w}}$$

$$E(\underline{w} + \delta \underline{w}) = E(\underline{w}) + (\delta \underline{w})^T \nabla E + \frac{1}{2} (\delta \underline{w})^T H (\delta \underline{w})$$

= 0, at the extremum

Extremum:

$$E(\underline{w} + \delta \underline{w}) = E(\underline{w}) + \frac{1}{2} (\delta \underline{w})^T H (\delta \underline{w})$$

→ a geometric interpretation