

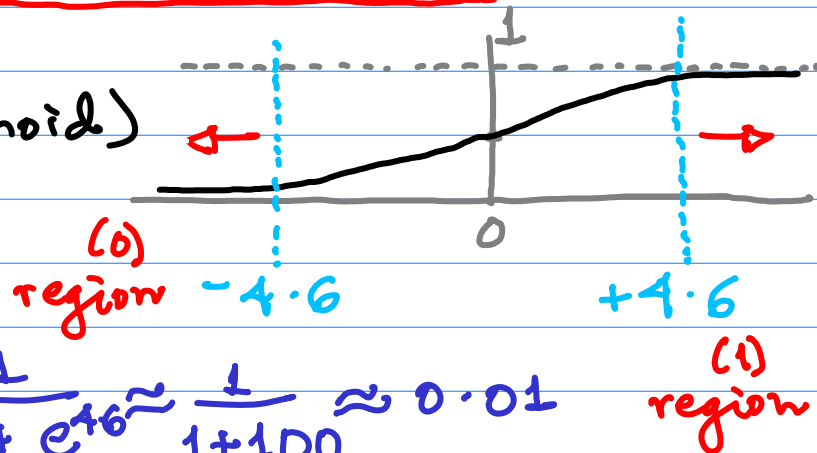
SOME HANDCRAFTED NN EXAMPLES

$$h(a) = \frac{1}{1 + e^{-a}} \quad (\text{sigmoid})$$

$$e^{4.6} \approx 100$$

When $a = -4.6$, $h(a) = \frac{1}{1 + e^{+4.6}} \approx \frac{1}{1 + 100} \approx 0.01$

When $a = +4.6$, $h(a) = \frac{1}{1 + e^{-4.6}} \approx \frac{1}{1 + 1/100} = \frac{100}{101} \approx 0.99$

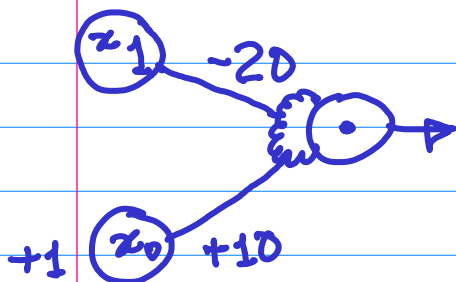


Example 1: NOT

$$a = \underline{w}^T \underline{x} + b$$

$$= -20x_1 + 10$$

x_1	\bar{x}_1
0	1
1	0



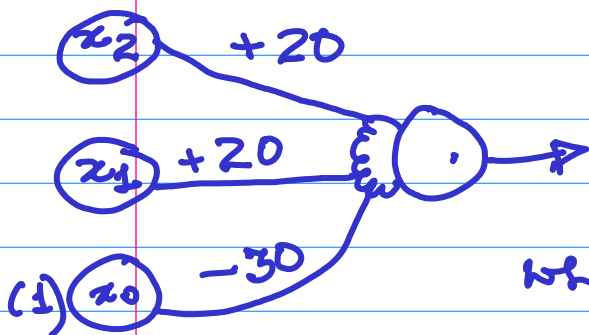
$$x_1 = 0, a = 10, h(+10) \approx 1$$

$$x_1 = 1, a = -10, h(-10) \approx 0$$

Example 2: AND

$$a = \underline{w}^T \underline{x} + b$$

x_2	x_1	$x_2 \cdot x_1$
0	0	0
0	1	0
1	0	0
1	1	1



$$a = 20x_2 + 20x_1 - 30$$

When $x_2 = 0, x_1 = 0$: $a = -30$

$$h(-30) \approx 0$$

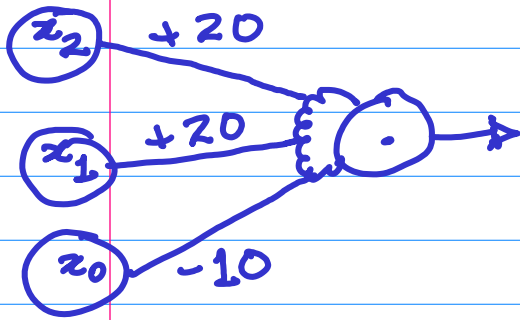
When $x_2 = 0, x_1 = 1$, $a = 20(0) + 20(1) - 30 = -10, h(-10) \approx 0$

When $x_2=1, x_1=0, a=20(1)+20(0)-30=-10, h(-10) \approx 0$

When $x_2=1, x_1=1, a=20(1)+20(1)-30=+10, h(+10) \approx 1$

Example 3: OR

x_2	x_1	$x_2 + x_1$
0	0	0
0	1	1
1	0	1
1	1	1



$$a = \underline{w}^T \underline{x} + b$$

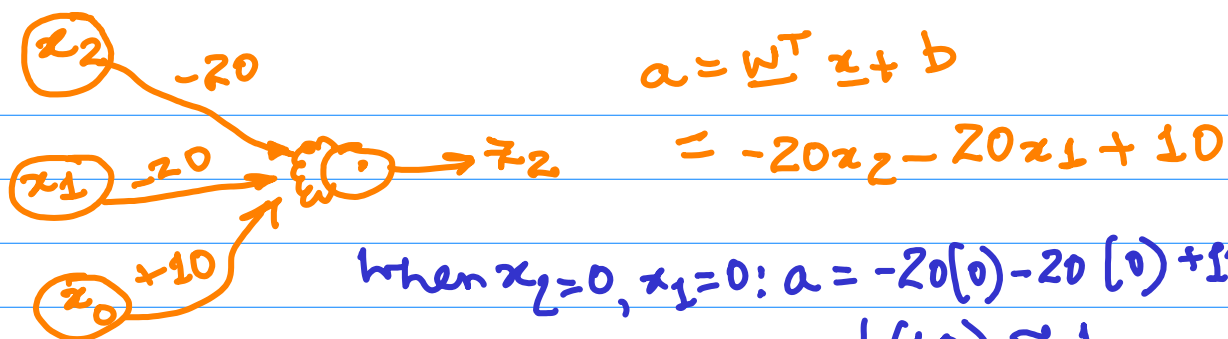
$$a = 20x_2 + 20x_1 - 10$$

When $x_2=0, x_1=0: a=20(0)+20(0)-10=-10, h(-10) \approx 0$

When $x_2=0, x_1=1: a=20(0)+20(1)-10=+10, h(+10) \approx 1$

When $x_2=1, x_1=0: a=20(1)+20(0)-10=+10, h(+10) \approx 1$

When $x_2=1, x_1=1: a=20(1)+20(1)-10=+30, h(+30) \approx 1$



When $x_2=0, x_1=0$: $a = -20(0) - 20(0) + 10 = 10$

$h(10) \approx 1$

When $x_2=0, x_1=1$: $a = -20(0) - 20(1) + 10 = -10$

$h(-10) \approx 0$

When $x_2=1, x_1=0$: $a = -20(1) - 20(0) + 10 = -10$

$h(-10) \approx 0$

When $x_2=1, x_1=1$: $a = -20(1) - 20(1) + 10 = -30$

$h(-30) \approx 0$

$z_2 = \overline{x_2} \cdot \overline{x_1}$

$y = \overline{x_2} \cdot \overline{x_1} + x_2 \cdot x_1 = x_2 \odot x_1$

(*) This is an informal (non-mathematical, intuitive) manifestation of a fundamental notion \rightarrow

"A feedforward NN with one hidden layer can represent any Boolean function"

(*) This also gives an intuitive (non-mathematical) perspective of a more general result: \rightarrow

"Multi-layer feedforward NNs with non-linear activation functions are universal approximators — they can approximate any function arbitrarily well."

Building Block (Input: 2-D: Image)

MNIST Numeral database: 28x28 images

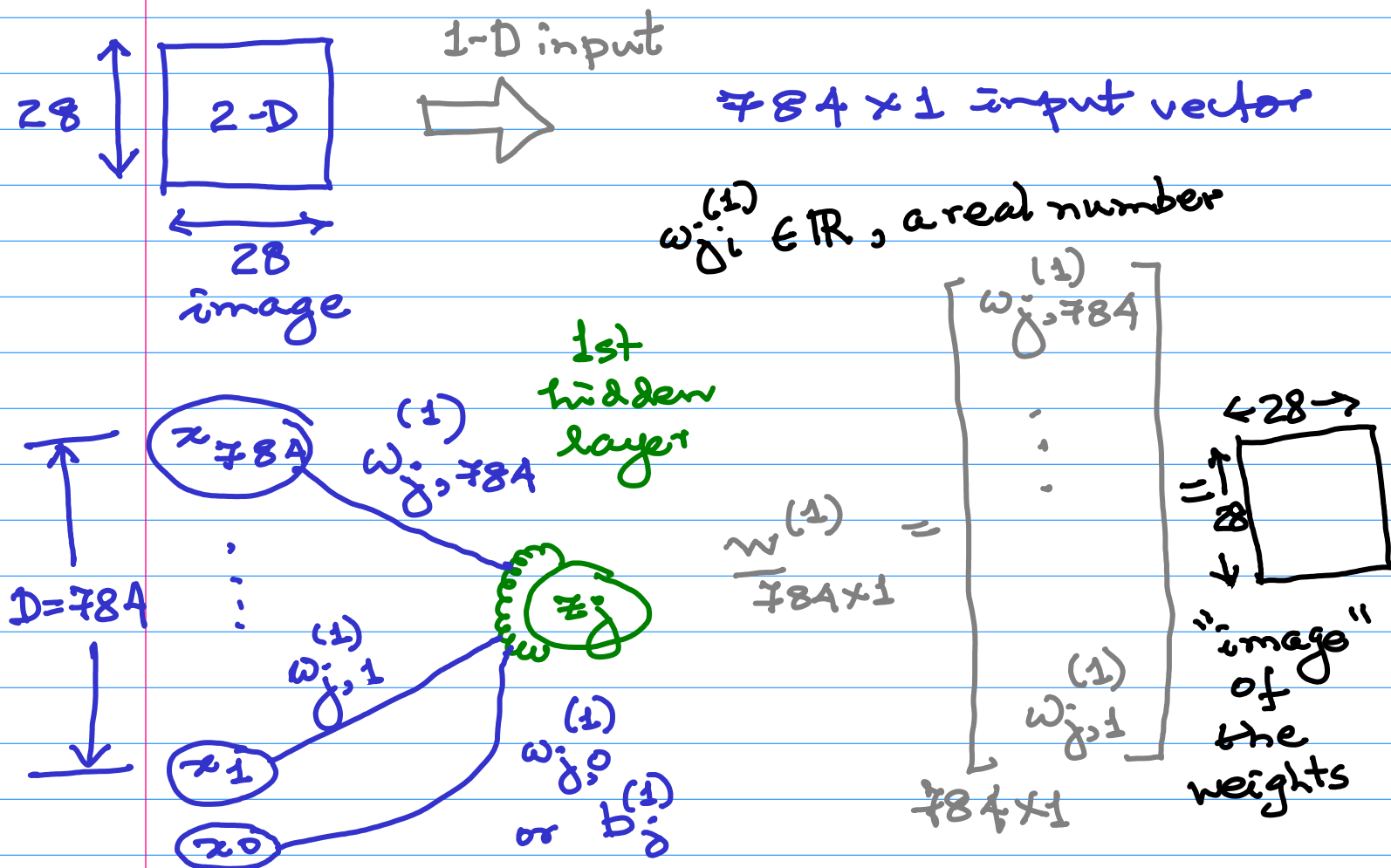
(grayscale: not binary (not 0 and 255, or normalised 0 and 1))
~ smooth

→ shades of grey as well, though most of the image is black or white.
(0) (255, or 1, normalised)

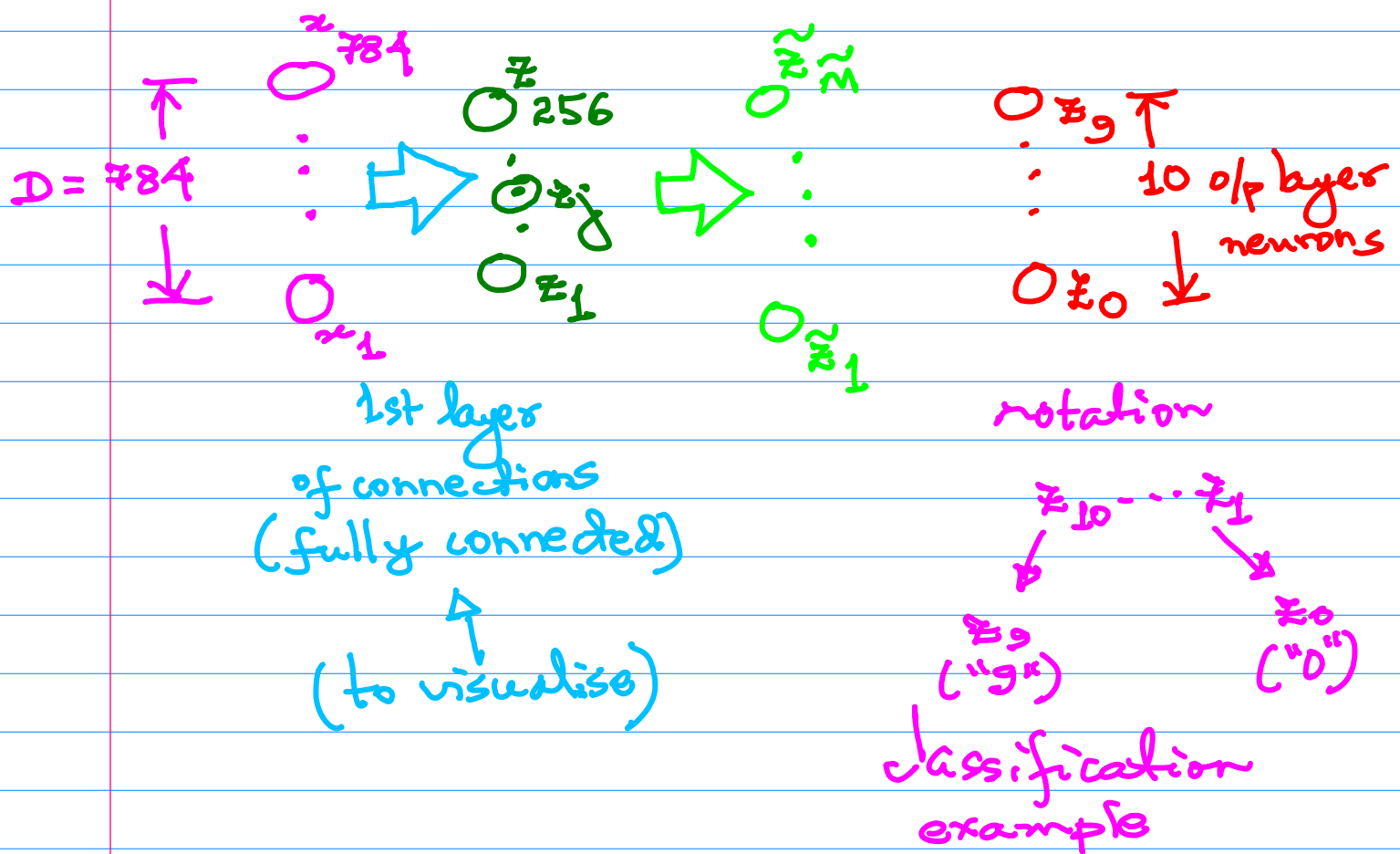
Images of the 10 numerals: 0 to 9

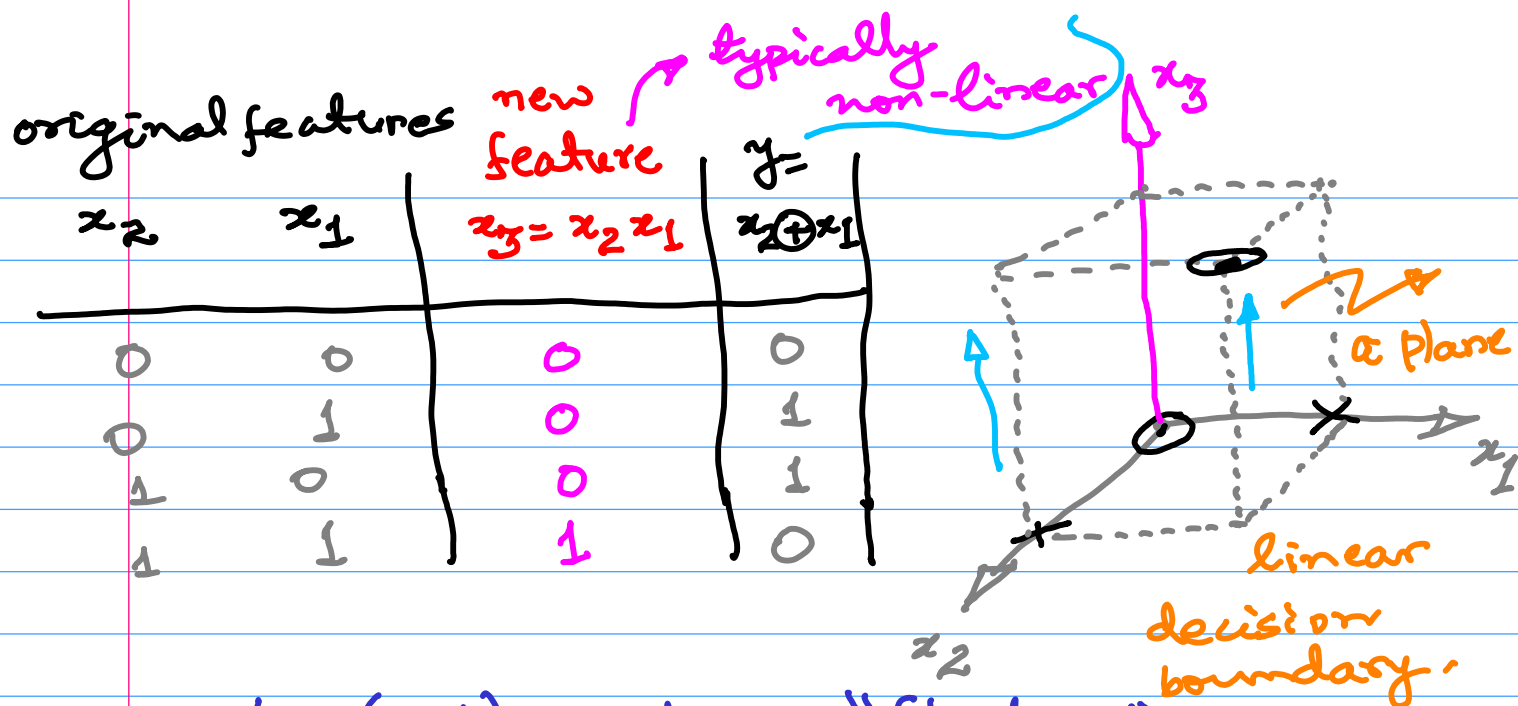
Basic structure: an MLP with 2 hidden layers

→ of specific interest is the first layer of connections



Take-home point: The input is not an ordinary 1-D vector, but a 2-D image
 \Rightarrow we can visualise the weights not as an ordinary 1-D vector, but a 2-D image with a similar spatial configuration / arrangement as the input itself.





The earlier (1,1) point now "floats up" leading to an infinite number of planes (linear decision boundaries in 3-D) now separating the two classes (much like the concentric circles doughnut 'floating up' over the vada)

The 'Factorisation' in Math/Summation

Where does this appear, and why? **Short Answer: Everywhere!**

Working Rule:-

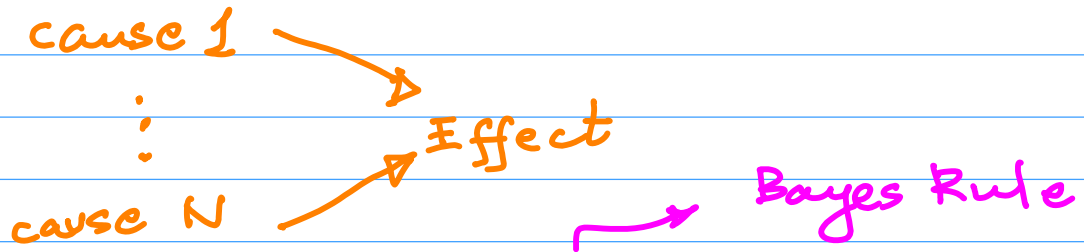
try putting it everywhere, and remove it if it is not required!

→ Probability (conditional/non-conditional)

→ Chain Rule (Calculus)
— total derivate/partial derivative.

PROBABILITY

The general probability factorisation



$$P(\text{cause \# } i | \text{effect}) = \frac{P(\text{effect} | \text{cause \# } i) P(\text{cause \# } i)}{P(\text{effect})}$$

Symmetrical

$$\underbrace{P(A|B) P(B)}_{P(A \text{ and } B)} = \underbrace{P(B|A) P(A)}_{P(B \text{ and } A)}$$

$$\boxed{P(A|B)} = \frac{P(B|A) \boxed{P(A)}}{P(B)}$$

a posteriori

probability of A

updated probability of A

$$x = x + 1$$

$$\Rightarrow x_{\text{new}} = x_{\text{old}} + 1$$

$$P(A)_{\text{updated}} = [\quad] \times P(A)_{\text{initial}} \rightarrow P(A)$$

$P(A|B)$

→ As such, there is no difference between a conditional probability and an unconditional probability → these are just the updated and initial variants of the same physical quantity.

$$p(A|B) = \frac{p(B|A) p(A)}{p(B)}$$

$$\boxed{p(A|B, C)} = \frac{p(B|A, C)}{p(B|C)} \boxed{p(A|C)}$$

updated prob of A initial prob of A

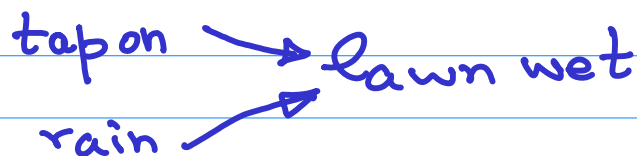
$$p(A)_{\text{updated}} \leftarrow \boxed{\text{Bayes Rule}} \leftarrow p(A)_{\text{initial}}$$

Where do the summations come in

$$p(\text{effect}) = \sum_j p(\text{effect} | \text{cause} \# j) p(\text{cause} \# j)$$

$\forall j$ → safely put a summation for all j 's

Example: →



Suppose we find the lawn to be wet in the morning.
 [Bayes Rule] $p(\text{rained last night} | \text{lawn wet in the morning})$

$$p(\text{rain}|\text{wet}) = \frac{p(\text{wet}|\text{rain}) p(\text{rain})}{p(\text{wet})}$$

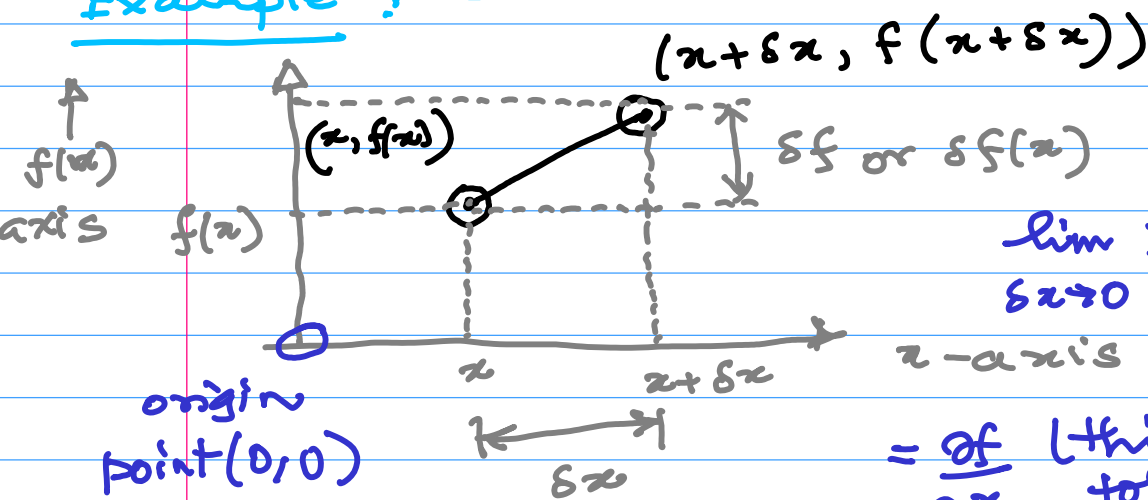
$$= \frac{p(\text{wet}|\text{rain}) p(\text{rain})}{p(\text{wet}|\text{rain}) p(\text{rain}) + p(\text{wet}|\text{tap on}) p(\text{tap on})}$$

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.

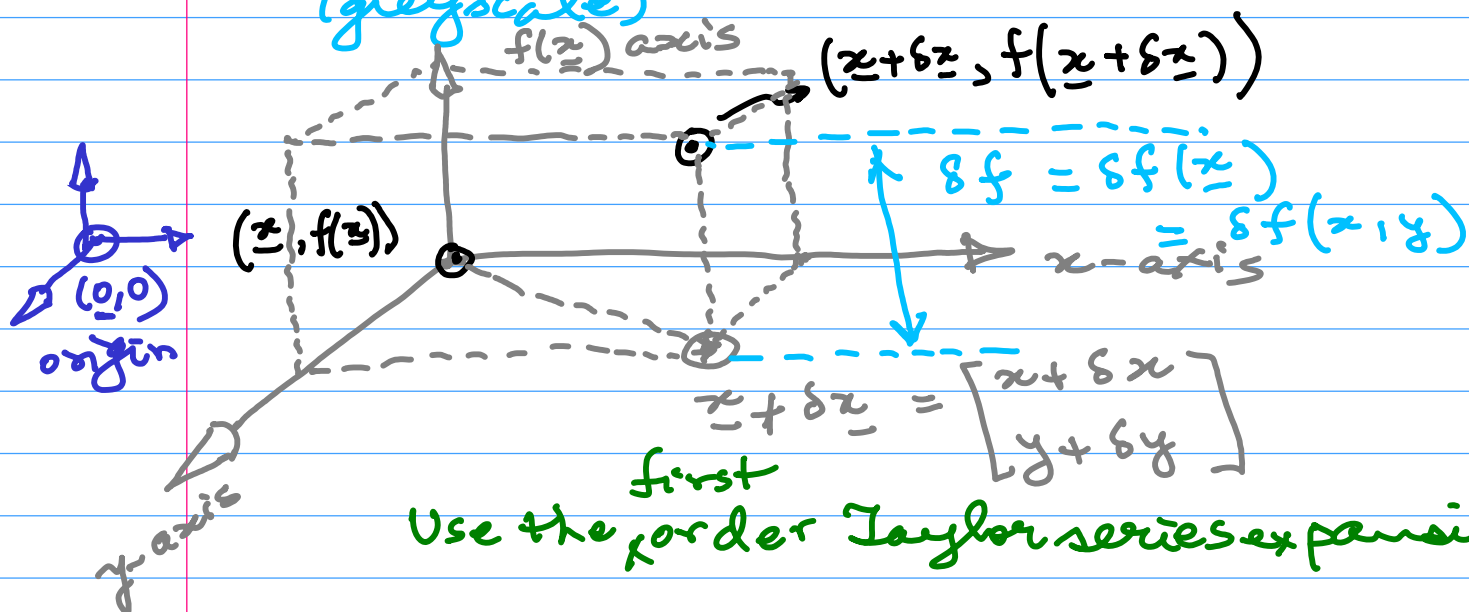
$$\lim_{\delta x \rightarrow 0} : \underbrace{f(x + \delta x) - f(x)}_{\delta f \text{ or } \delta f(x)} = \left(\frac{\partial f}{\partial x} \right) \underbrace{\delta x}_{\text{small change in the independent variable.}}$$

Small change in the dependent variable, or the function

[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 $\underline{x} = \begin{bmatrix} x \\ y \end{bmatrix}$
(grayscale)



Use the ^{first} order Taylor series expansion

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

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

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

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

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