

X-OR problem (contd.) Formulate this as a regression problem and use MSE loss (Mean Square Error)

Training Set $\underline{x} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}$ $X = [\underline{x}_{(1)} \quad \underline{x}_{(2)} \quad \underline{x}_{(3)} \quad \underline{x}_{(4)}] = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$

$\underline{t} = [t_{(1)} \quad t_{(2)} \quad t_{(3)} \quad t_{(4)}] = [0 \quad 1 \quad 1 \quad 0]$

MSE loss $J(\underline{w}) \triangleq \frac{1}{4} \sum_{n=1}^4 [\gamma(\underline{x}_{(n)}, \underline{w}) - t_{(n)}]^2$
↗ parameters

linear model

$y(\underline{x}) = \gamma(\underline{x}, \underline{w}) = \underbrace{\underline{w}^T \underline{x}}_{\text{homogeneous}} = \underbrace{\underline{w}^T \underline{x} + b}_{\text{non-homogeneous}} \text{ or } \omega_0$

(*) $\frac{\partial J(\underline{w})}{\partial \underline{w}} = \frac{1}{4} \cdot 2 \sum_{(n)} \{ \gamma(\underline{x}_{(n)}, \underline{w}) - t_{(n)} \} \frac{\partial \gamma(\underline{x}_{(n)}, \underline{w})}{\partial \underline{w}} = 0$

$\Rightarrow \frac{\partial (\underline{w}^T \underline{x})}{\partial \underline{w}} = \underline{x} = \frac{\partial (\underline{x}^T \underline{w})}{\partial \underline{w}} \quad (\text{Rule})$

$\Rightarrow \boxed{\frac{1}{2} \sum_{(n)} \{ \gamma(\underline{x}_{(n)}, \underline{w}) - t_{(n)} \} \underline{x}_{(n)} = 0} \quad \text{--- (1)}$

(*) $\frac{\partial J(\underline{w})}{\partial b} = \frac{1}{4} \cdot 2 \sum_{(n)} \{ \gamma(\underline{x}_{(n)}, \underline{w}) - t_{(n)} \} \cdot 1 = 0$
 $\frac{\partial (\underline{w}^T \underline{x} + b)}{\partial b} = 1$

$\Rightarrow \boxed{\sum_{(n)} \{ \underline{w}^T \underline{x}_{(n)} + b - t_{(n)} \} = 0} \quad \text{--- (2)}$

$$\begin{aligned}
 (1) \Rightarrow & \{w_2(0) + w_1(0) + b - 0\} \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\
 & + \{w_2(0) + w_1(1) + b - 1\} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\
 & + \{w_2(1) + w_1(0) + b - 1\} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\
 & + \{w_2(1) + w_1(1) + b - 0\} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}
 \end{aligned}$$

$$\Rightarrow \begin{bmatrix} w_2 + b - 1 + w_2 + w_1 + b \\ w_1 + b - 1 + w_2 + w_1 + b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 2w_2 + w_1 + 2b \\ 2w_1 + w_2 + 2b \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{cases} 2w_2 + w_1 + 2b = 1 \\ 2w_1 + w_2 + 2b = 1 \end{cases} \quad (1)$$

$$\begin{aligned}
 (2) \Rightarrow & w_2(0) + w_1(0) + b - 0 \\
 & + w_2(0) + w_1(1) + b - 1 \\
 & + w_2(1) + w_1(0) + b - 1 \\
 & + w_2(1) + w_1(1) + b - 0 = 0
 \end{aligned}$$

$$\Rightarrow 4b + 2(w_2 + w_1) = 2 \Rightarrow 2b + (w_2 + w_1) = 1 \quad (2)$$

Put (2) in (1) to get

$$\begin{aligned}
 & \left[\begin{array}{l} 2w_2 + \cancel{w_1} + \cancel{1} - (w_2 + \cancel{w_1}) = \cancel{1} \\ 2w_1 + \cancel{w_2} + \cancel{1} - (\cancel{w_2} + w_1) = \cancel{1} \end{array} \right] \Rightarrow \begin{aligned} w_2 &= 0 \\ w_1 &= 0 \\ b &= 0.5 \end{aligned}
 \end{aligned}$$

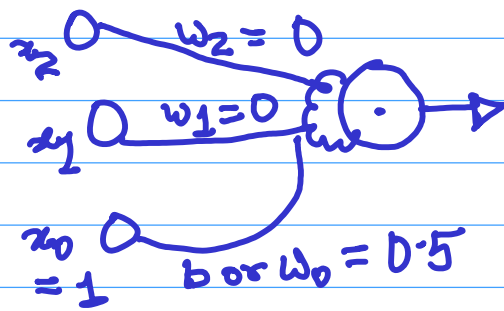
$$\underline{w} = \underline{0}, \quad b = 0.5$$

$$y(n) = \underline{w}^T \underline{x}(n) + 0.5$$

whenever is \underline{x}
the output should be 0.5

What did we try?

$y = \sigma(a) = a$ unit function



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

(non-homogeneous)

$$= w_2 x_2 + w_1 x_1 + b$$

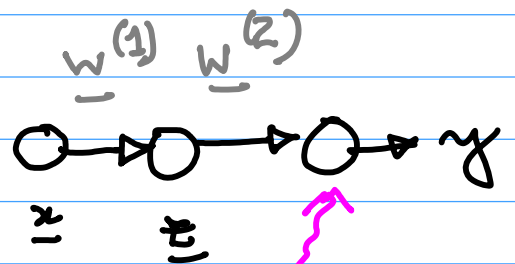
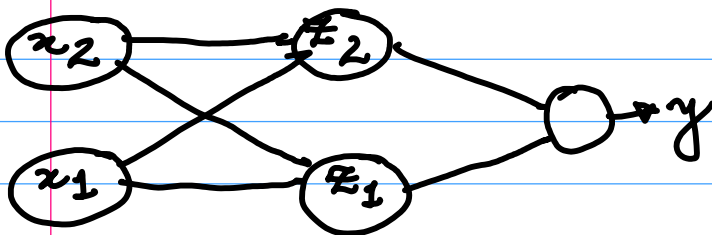
$$= 0.5 \text{ always}$$

Some points: →

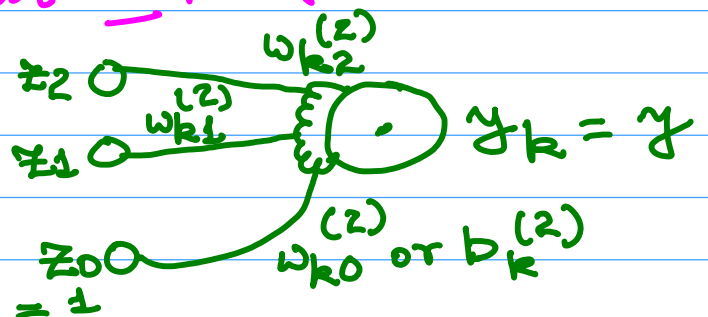
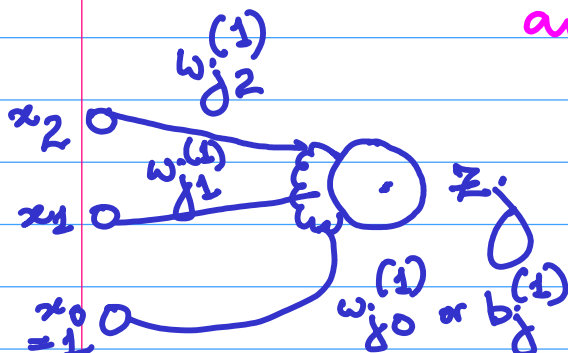
- This is a simple example with regression (though the problem is actually one of classification)
- This is one of the myriad possible solutions
- Unnecessarily trying to build a classifier out of a regressor (that too, for binary inputs)

- Handcrafted

ONE possible solution: -



regressor applied to $\underline{z} = \phi(\underline{x})$,
and not \underline{x} itself



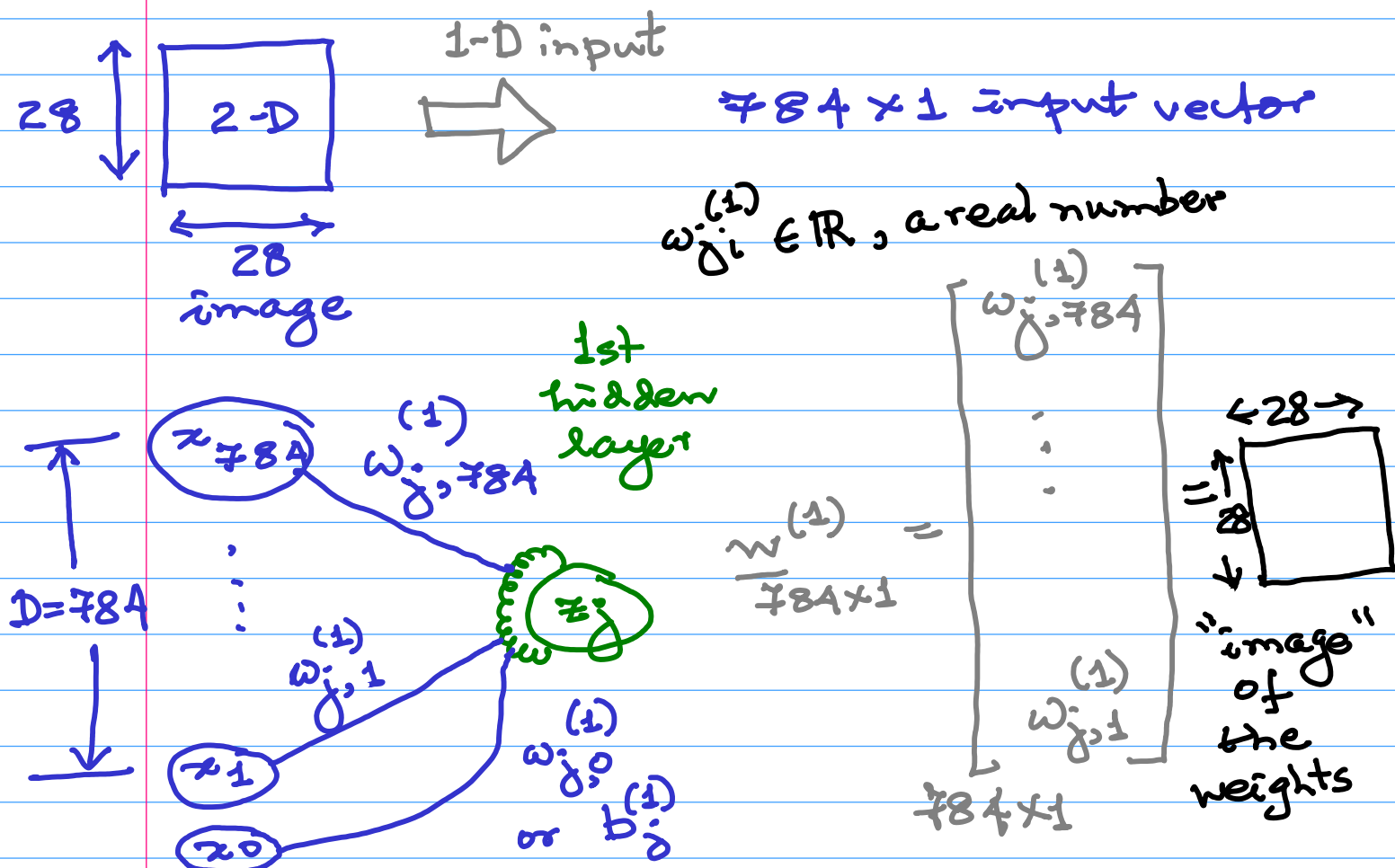
Building Block (Input: 2-D: Image)

MNIST Numeral database: 28×28 images
(grayscale: not binary (not 0 and 255, or normalised 0 and 1))
~ smooth

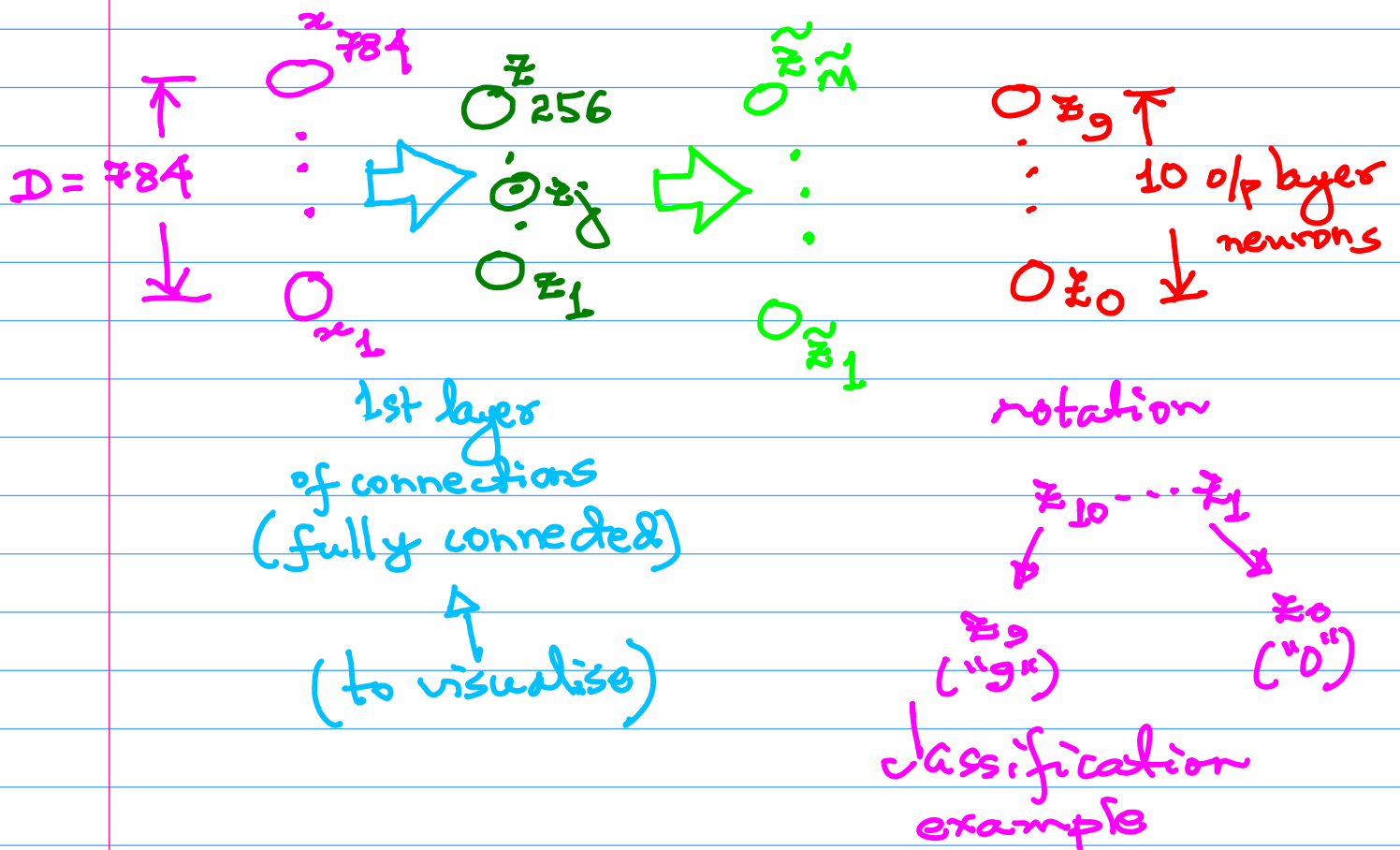
→ shades of grey as well, though most of the image is black or white.
(0) (255, or 1, normalised)

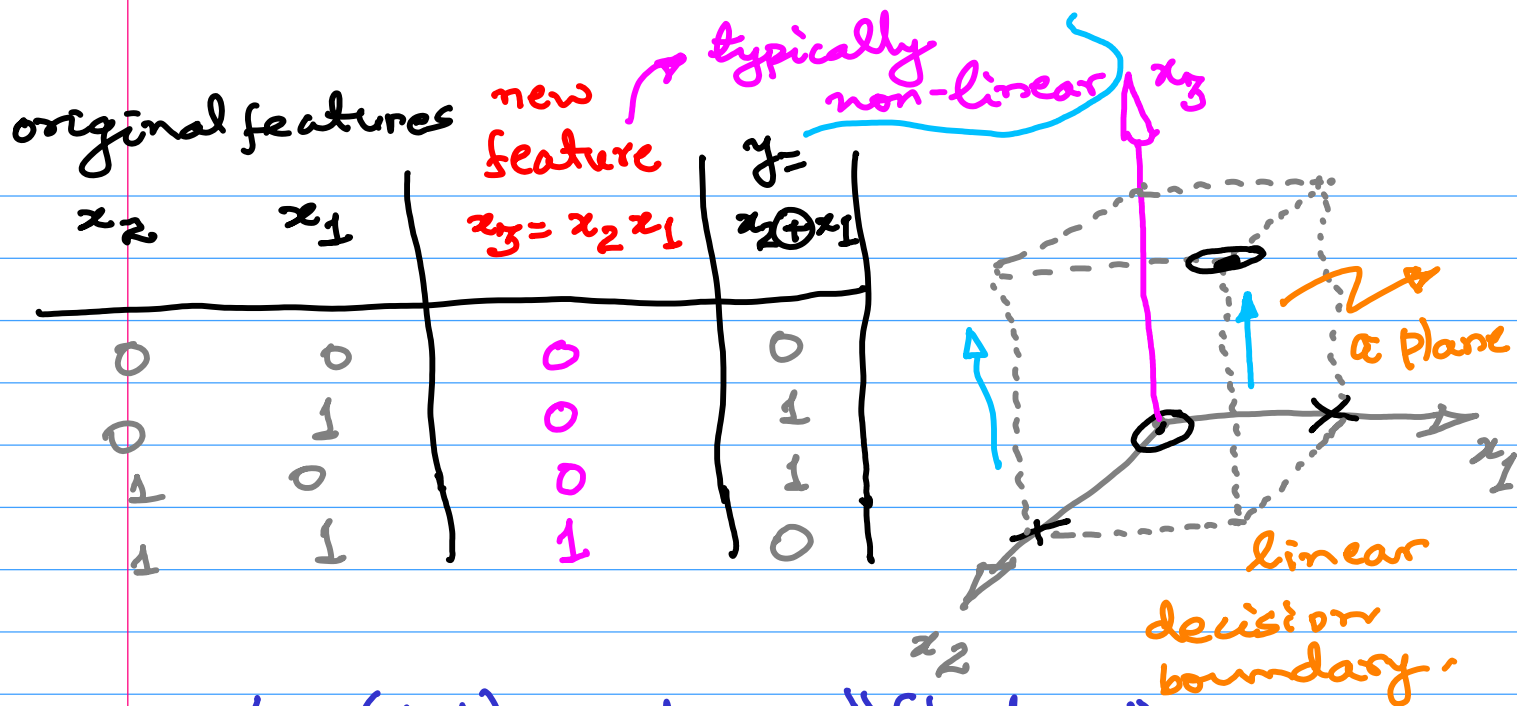
Images of the 10 numerals: 0 to 9

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:-

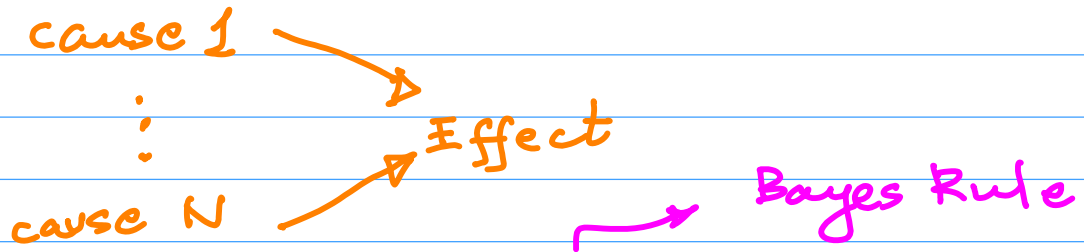
try putting it everywhere, and remove it if it is not required!

→ 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

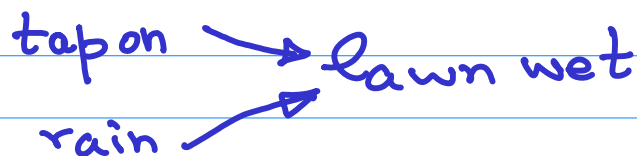
$$p(A)_{\text{updated}} \leftarrow \boxed{\text{Bayes Rule}} \leftarrow p(A)_{\text{initial}}$$

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] $p(\text{rained last night} | \text{lawn wet in the morning})$

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

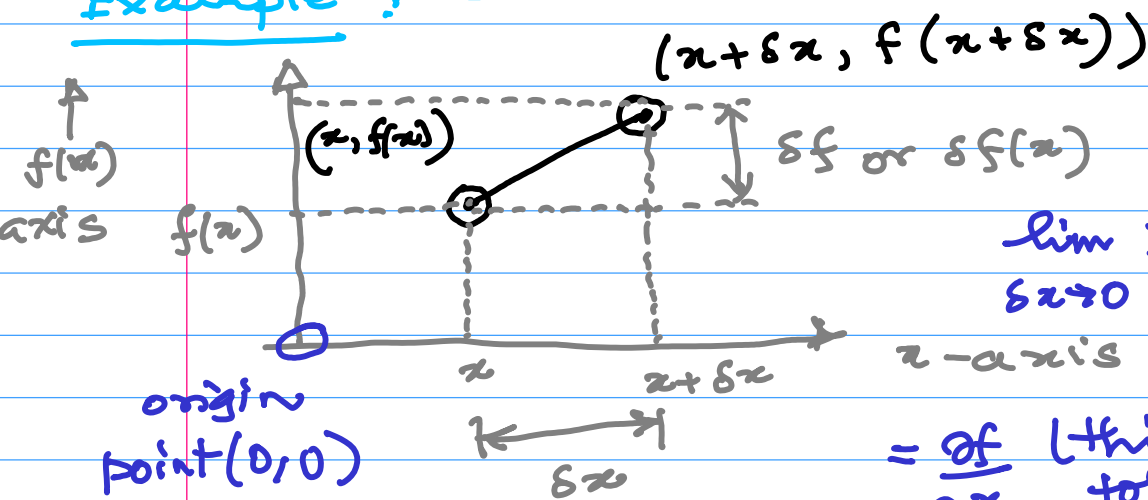
$$= \frac{p(\text{wet}|\text{rain}) p(\text{rain})}{p(\text{wet}|\text{rain}) p(\text{rain}) + p(\text{wet}|\text{tap on}) p(\text{tap on})}$$

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 $\underline{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

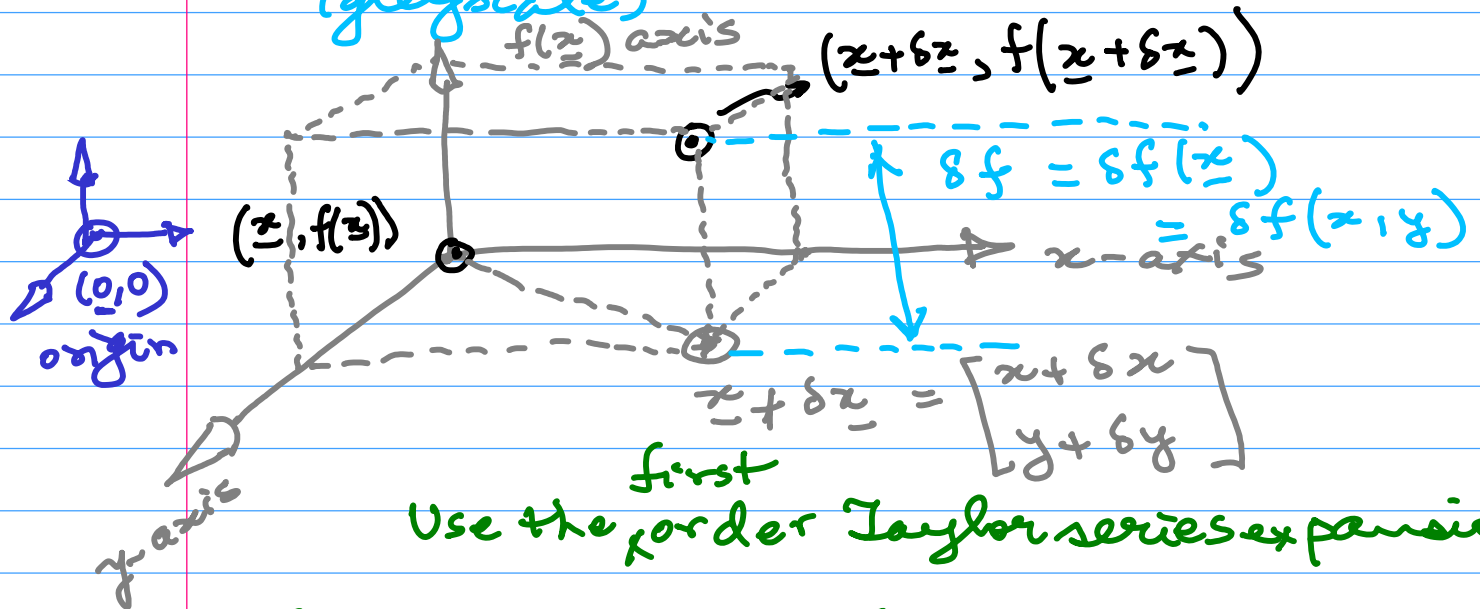
$\underline{x} \rightarrow$ vector

one dependent variable $f(\underline{x})$

$f(\underline{x}) \rightarrow$ 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)



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

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

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

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

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

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

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

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

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

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

Generalise to a function of D variables

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

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

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

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

BACK PROPAGATION

$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - \eta \nabla E$$

1849, Cauchy

Error function
What is this?

Gradient Descent

We typically start with random weights

$$\begin{bmatrix} \omega_{ji}^{(1)} \\ \vdots \\ \omega_{kj}^{(2)} \\ \vdots \end{bmatrix}$$

Training:
(Supervised)

[training] inputs \rightarrow outputs
To learn the weights \equiv function
which the NN implements

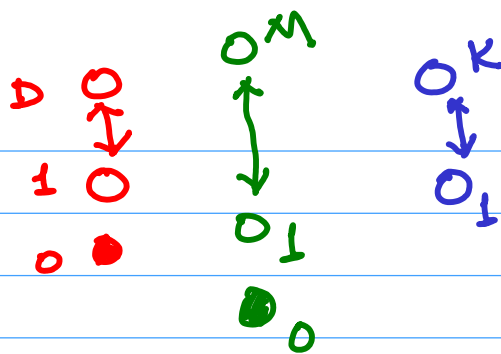
Error: b/w the ideal output
& the one which the NN
gives us currently.

N points
index $n = 1 \dots N$

$$E(\underline{w}) = \sum_{n=1}^N E_n(\underline{w})$$

$$\frac{\partial E(\underline{w})}{\partial \omega} = \sum_n \frac{\partial E_n}{\partial \omega}$$

Example:



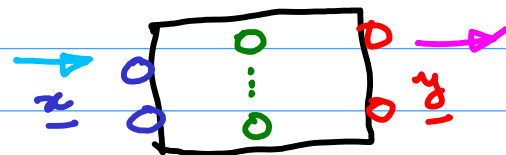
Assumptions: sum of squares errors

- Hidden layer activation function $h(a) = \tanh(a)$
- Output layer activation function $\sigma(a) = a, \forall k = a_k$

$$h(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

$$\begin{aligned} h'(a) &= \frac{\partial h(a)}{\partial a} = \frac{(e^a + e^{-a})(e^a - (-e^{-a})) - (e^a - e^{-a})(e^a - e^{-a})}{(e^a + e^{-a})^2} \\ &= \frac{(e^a + e^{-a})^2 - (e^a - e^{-a})^2}{(e^a + e^{-a})^2} = 1 - \left[\frac{e^a - e^{-a}}{e^a + e^{-a}} \right]^2 = 1 - h^2(a) \end{aligned}$$

For the n 'th training data item:



$$E_n \triangleq \frac{1}{2} \sum_{k=1}^K (\gamma_k - t_k)^2$$

→ ground truth / label

$\gamma_k = k$ 'th item of the function $\gamma(\underline{w}, \underline{x})$

$t_k = k$ 'th item of the target vector

why only γ_k ?
(defined only at the output layer)