

Intuitive Explanation for the weight update eqn:

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

[Bauchy, 1849]

$\underline{w} \rightarrow$ all weights of a neural network

1 layer: perceptron

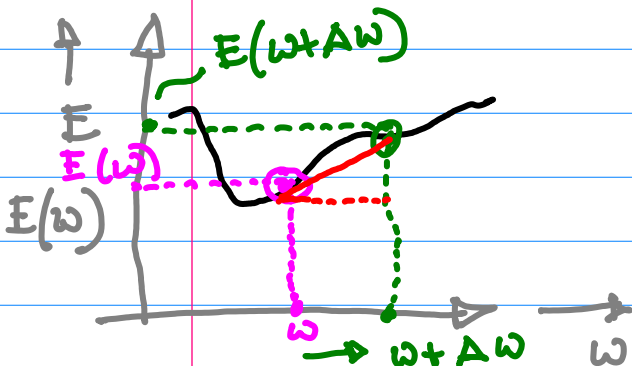
many layers: MLP

e.g., $w_{ji}^{(1)}$ and $w_{kj}^{(2)}$

∇E scalar \rightarrow vector of all partial derivatives $\left[\frac{\partial E}{\partial w} \right]$ 'gradient'

$E(\underline{w})$ is an error function of all weights \underline{w}

\rightarrow This equation is DIMENSIONALLY consistent



$$\text{slope} = \frac{E(w + \Delta w) - E(w)}{(w + \Delta w) - (w)}$$

$$= \frac{E(w + \Delta w) - E(w)}{\Delta w}$$

$$\lim_{\Delta w \rightarrow 0} \rightarrow \frac{\partial E}{\partial w}$$

Intuition: To get to a local min of the error function, we must go **against** the gradient / the slope \rightarrow intuition behind the '-' sign

What is η ? ('step size')

Small $\eta \rightarrow$ small steps, longer to get to the min

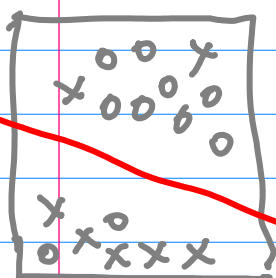
large $\eta \rightarrow$ large steps, we may miss the local min

η : step size / learning rate.

A limitation of the Perceptron: seeks a linear decision boundary

$$y(\underline{x}, \underline{w}) = \text{sgn} \left(\sum_{i=1}^D w_i x_i + b \right)$$

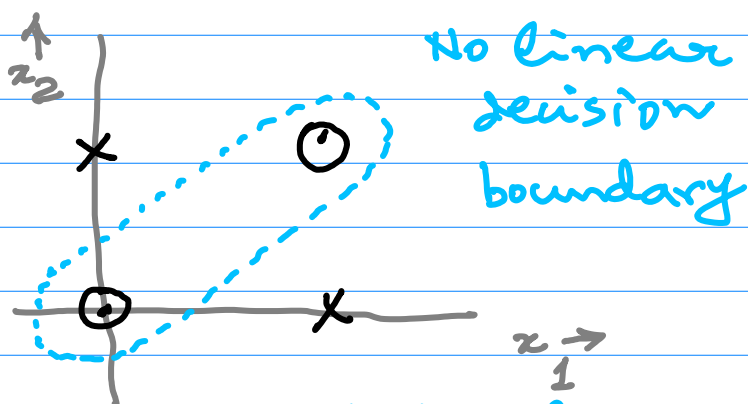
$$= \text{sgn} \left(\underbrace{\underline{w}^T \underline{x} + b}_{\text{linear eqn}} \right)$$



↪ The linear decision boundary is a limitation

'The XOR problem' (Marvin Minsky)

x_2	x_1	$y = x_2 \oplus x_1$
0	0	0 → ○
0	1	1 → ×
1	0	1 → ×
1	1	0 → ○



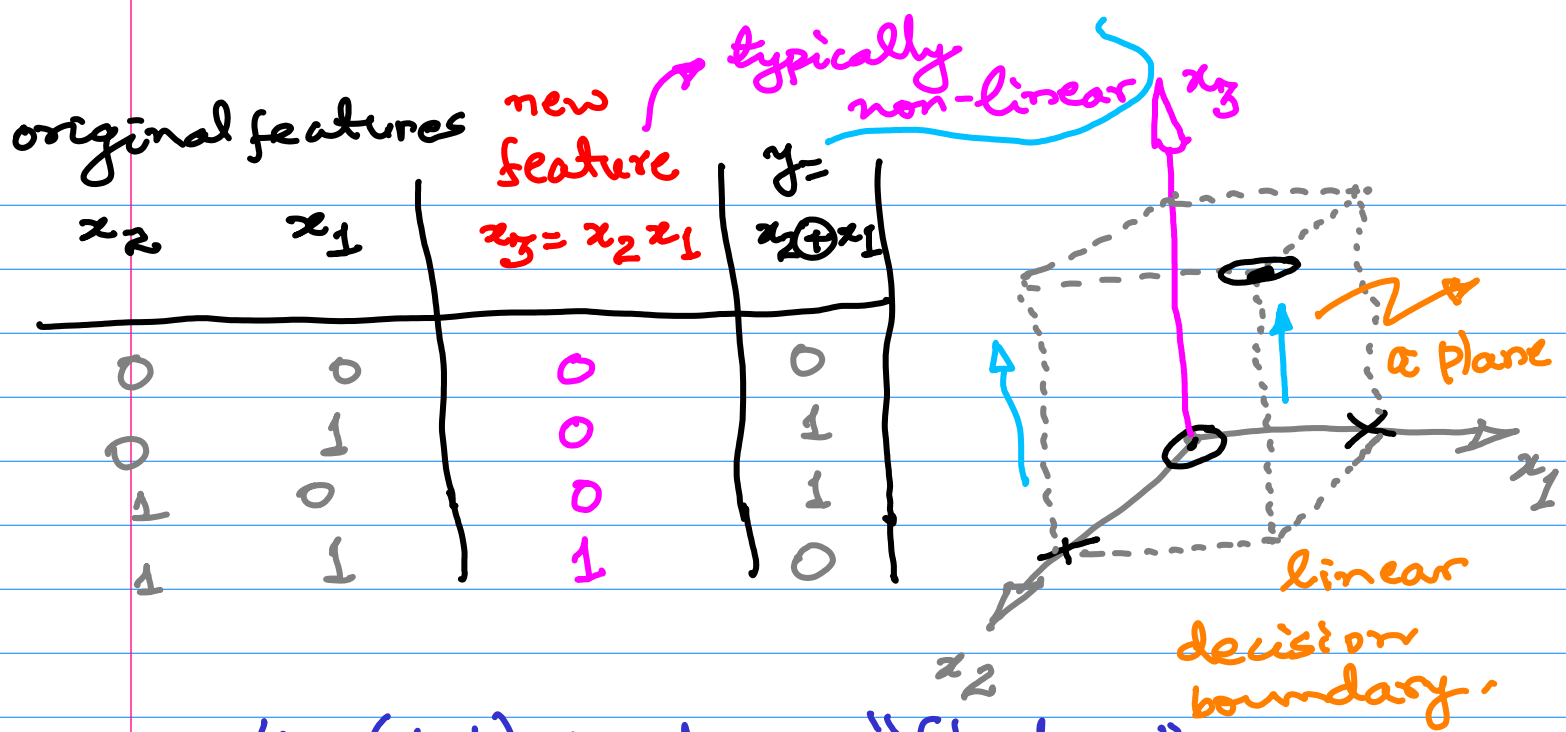
separable, but not linearly separable.

⇒ cannot be solved with a linear classifier such as a basic perceptron.

However, it can be solved by

- Feature transformation + perceptron
- by hand!
- MLP (hidden layer(s))

this itself has a feature transformation connotation



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!

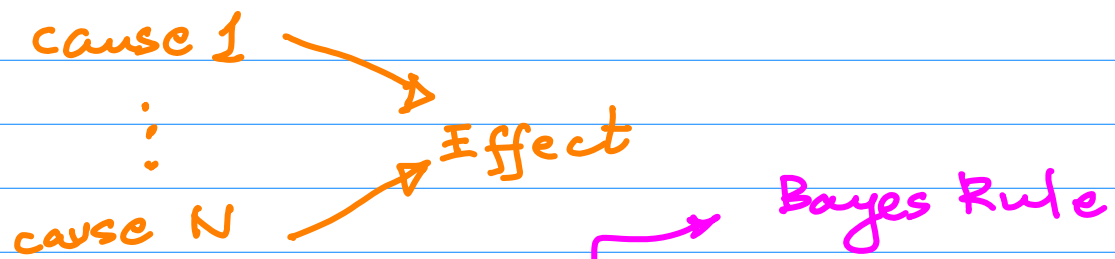
→ Probability (conditional/non-conditional)

→ 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

a priori 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

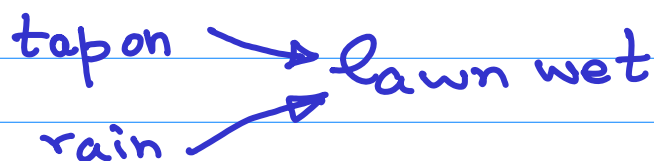
$$p(A)_{\text{updated}} \leftarrow \boxed{\text{Bayes Rule}} \leftarrow p(A)_{\text{initial}}$$

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

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

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