original features | new feature (typically non-linear)

| $x_2$ | $x_1$ | $z_3 = x_2 x_1$ | $y = x_2 \oplus x_1$ |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 |

$x_3$

a plane

linear decision boundary.

$x_1$

$x_2$

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!**
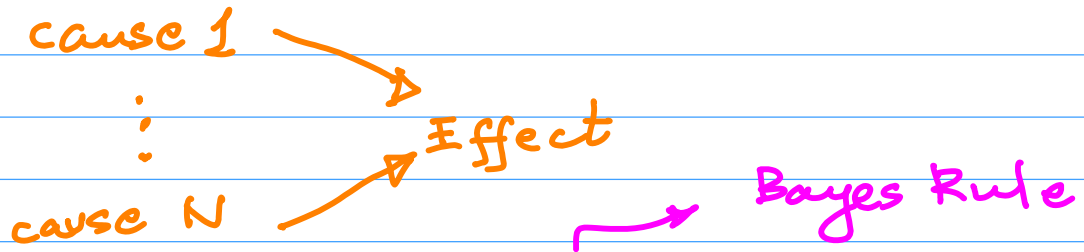
Working Rule :-
try putting it everywhere, and remove it if it is not required!

→ Probability (conditional / non-conditional)

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

# PROBABILITY

## The general probability factorisation

cause 1

⋮

cause N → Effect

Bayes Rule

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

Symmetrical

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

$$\boxed{P(A|B)} = \frac{P(B|A) \boxed{P(A)}}{P(B)} \longrightarrow \text{a priori probability of } A$$

a posteriori probability of A

updated probability of A

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

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

$$P(A|B)$$

→ 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.

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

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

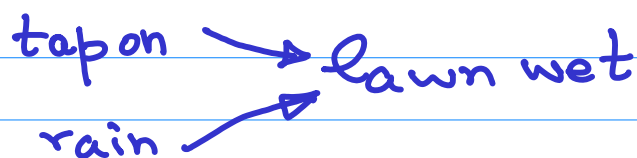updated prob of A

initial prob of A

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

where do the summations come in

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

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

Example! →

tap on → lawn wet

rain →

Suppose we find the lawn to be wet in the morning.  $p\left(\text{rained last night} \,\middle|\, \text{lawn wet in the morning}\right)$

[Bayes Rule]

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

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

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

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

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

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

$$= \sum \left(\frac{\partial f}{\partial \text{var}}\right)(\delta\text{var})$$

$$\text{var} = x, y, z$$

## 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 $\delta t \to 0$

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

Here, we had <u>one</u> variable $t$, so the partial derivative is also the total derivative.

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

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

(all the $x_i$'s are functions of $t_j$

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

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

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

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

# 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

$$E(\underline{w} + \delta\underline{w}) = E(\underline{w}) + (\overline{\nabla}E) \cdot (\delta\underline{w})$$

error function
minimise

weights
(All)

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

→ vector

scalar

At the extremum, $\overline{\nabla}E = \underline{0}$

max/min/saddle
point vector

$$\begin{bmatrix} \partial E/\partial w_2 \\ \partial E/\partial w_1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

(2-D)

## Local Quadratic Approximation

higher order terms

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

$$(\delta\underline{w})^T \bar{\nabla}E$$

↳ 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

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

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

$$= \underline{0}, \text{ at the extremum}$$

Extremum:

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

→ a geometric interpretation

# BACK PROPAGATION

$$w^{(\tau+1)} = w^{(\tau)} - \eta \, \nabla E$$

1849, Cauchy

Error function
What is this?

$$w_{ji}^{(1)}$$

$$w_{kj}^{(2)}$$

Training: (supervised)

[training] inputs → outputs

To learn the weights ≡ function which the NN implements

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

**Gradient Descent**

We typically start with random weights

N points
index $n = 1 .. N$

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

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

**Example:**

$O^D$   $\uparrow\downarrow$ $O$
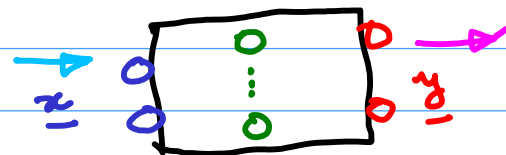 1  $O$
 0  ●

$O^M$ ↑↓ $O_1$ ● $O$

$O^K$ ↑↓ $O_1$

**Assumptions:** sum-of-squares errors

→ Hidden layer activation function
$$h(a) = \tanh(a)$$

→ Output layer activation function
$$\sigma(a) = a, \quad y_k = a_k$$

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

$$h'(a) = \frac{\partial h(a)}{\partial a} = \frac{(e^a + e^{-a})(e^a - (-e^{-a})) - (e^a - e^{-a})(e^a - e^{-a})}{(e^a + e^{-a})^2}$$

$$= \frac{(e^a + e^{-a})^2 - (e^a - e^{-a})^2}{(e^a + e^{-a})^2} = 1 - \left[\frac{e^a - e^{-a}}{e^a + e^{-a}}\right]^2 = 1 - h^2(a)$$

**For the n'th training data item:**



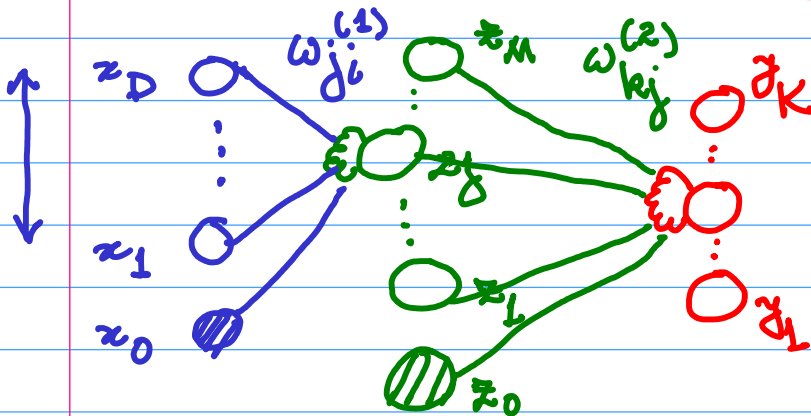$$E_n \triangleq \frac{1}{2}\sum_{k=1}^{k}(y_k - t_k)^2$$

→ ground truth / label

↳ why only $y_k$?
(defined only at the output layer)

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

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

# BACKPROPAGATION (contd.) [An example]



[2 layers of connections]

## Why is this called 'Backpropagation'?



Assume
→ $h(\cdot)$: activation function is $\tanh()$
→ $\sigma(\cdot)$: activation function is a unit function

Training phase: error is at the output: this is fed back to update the weights

$n \in \{1, N\}$

FOR each pattern in the training set in turn, we follow the following operations.

① A 'forward propagation'

($j$: hidden layer)

$$a_j^{(1)} = \underline{w}_j^{(1)T} \underline{x} = \sum_{i=0}^{D} w_{ji}^{(1)} x_i$$

$$z_j = \tanh(a_j)$$

$$y_k = \underline{w}_k^{(2)T} \underline{z} = \sum_{j=0}^{M} w_{kj}^{(2)} z_j$$

② Evaluate $\delta_k$'s for each output unit

$$\delta_k = y_k - t_k \quad \text{What is this, and how?}$$

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

$$\delta_k \triangleq \frac{\partial E_n}{\partial a_k} = \frac{1}{2} \cdot 2 (y_k - t_k) \underbrace{\frac{\partial y_k}{\partial a_k}}_{} \rightarrow = 1 \text{ as}$$

output layer
activation

$y_k = a_k$

$(\sigma(\cdot) = \text{unit fn})$

$$\delta_k = y_k - t_k \quad (\text{Else, according to the specific activation function } \sigma(\cdot) \text{ at the output layer})$$

③ <u>Backpropagate these to obtain $\delta_j$'s for the hidden layer units</u>

$$\boxed{\delta_j = (1 - z_j^2) \sum_{k=1}^{K} w_{kj} \cdot \boxed{\delta_k}}$$

What is this, and how?

$\rightarrow$ previous step (step #2) [output layer]

$\rightarrow$ [hidden layer]

○
⋮    ⟸    ○
○         ⋮
              ○

$$a_k = \underline{w}^{(2)T} \underline{z}$$
$$= \sum_{j=0}^{M} w_{kj}^{(2)} z_j$$
$$= \sum_{j=0}^{M} w_{kj}^{(2)} h(a_j)$$

$$\delta_j \triangleq \frac{\partial E_n}{\partial a_j} = \sum_{\forall k} \boxed{\frac{\partial E_n}{\partial a_k}} \boxed{\frac{\partial a_k}{\partial a_j}}$$

$\delta_k$

[step #2]

$$\Rightarrow \frac{\partial a_k}{\partial a_j} = w_{kj}^{(2)} \frac{\partial h(a_j)}{\partial a_j} = w_{kj}^{(2)} h'(a_j)$$
$$= w_{kj}^{(2)} (1 - z_j^2)$$

$$\delta_j = \sum_{\forall k} \delta_k w_{kj}^{(2)} (1 - z_j^2) = (1 - z_j^2) \sum_{\forall k} w_{kj}^{(2)} \delta_k$$