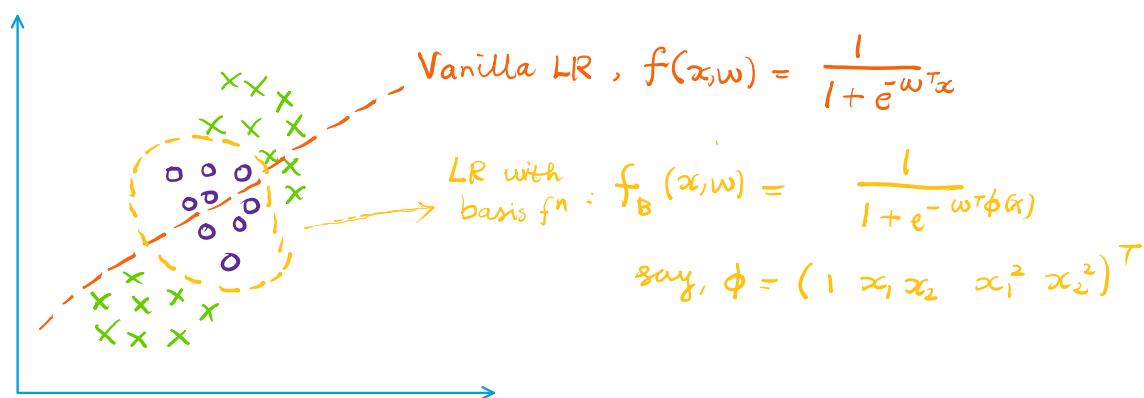
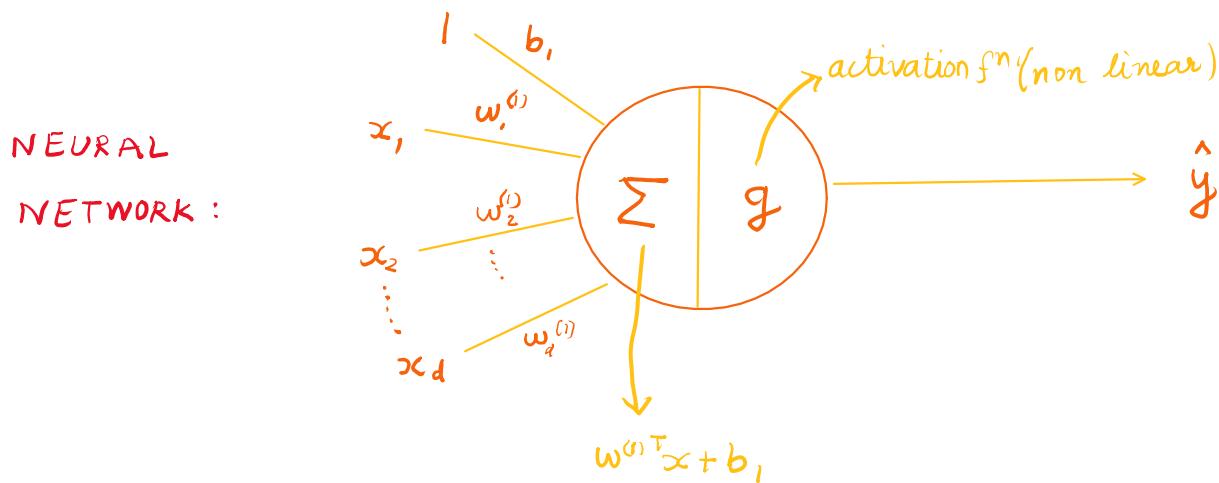
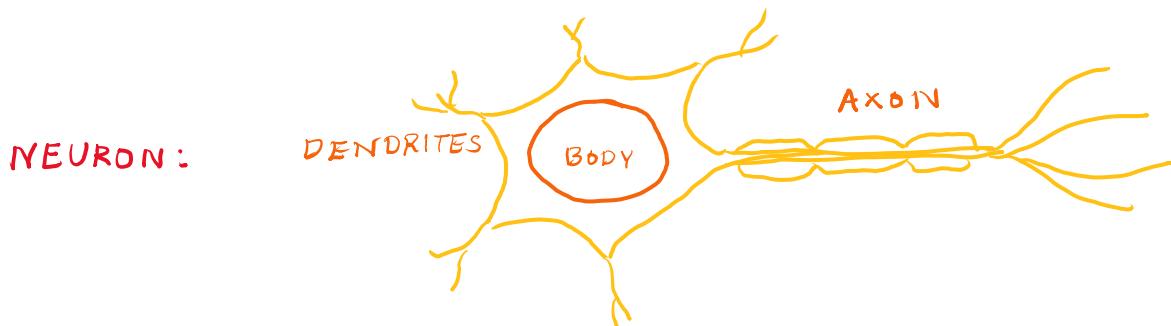


Example 1:



The underlying philosophy of Neural Networks is to come up with some basis fn  $\phi(x)$  without explicitly programming it.

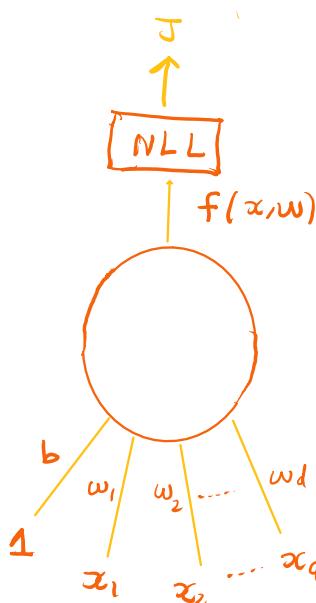
Origin of NN: Started in 1943 then progressed during 60's & 80's.  
NNs are inspired by Neurons.



Logistic regression using NN

$$f(x, w) = g(w^\top x + b)$$

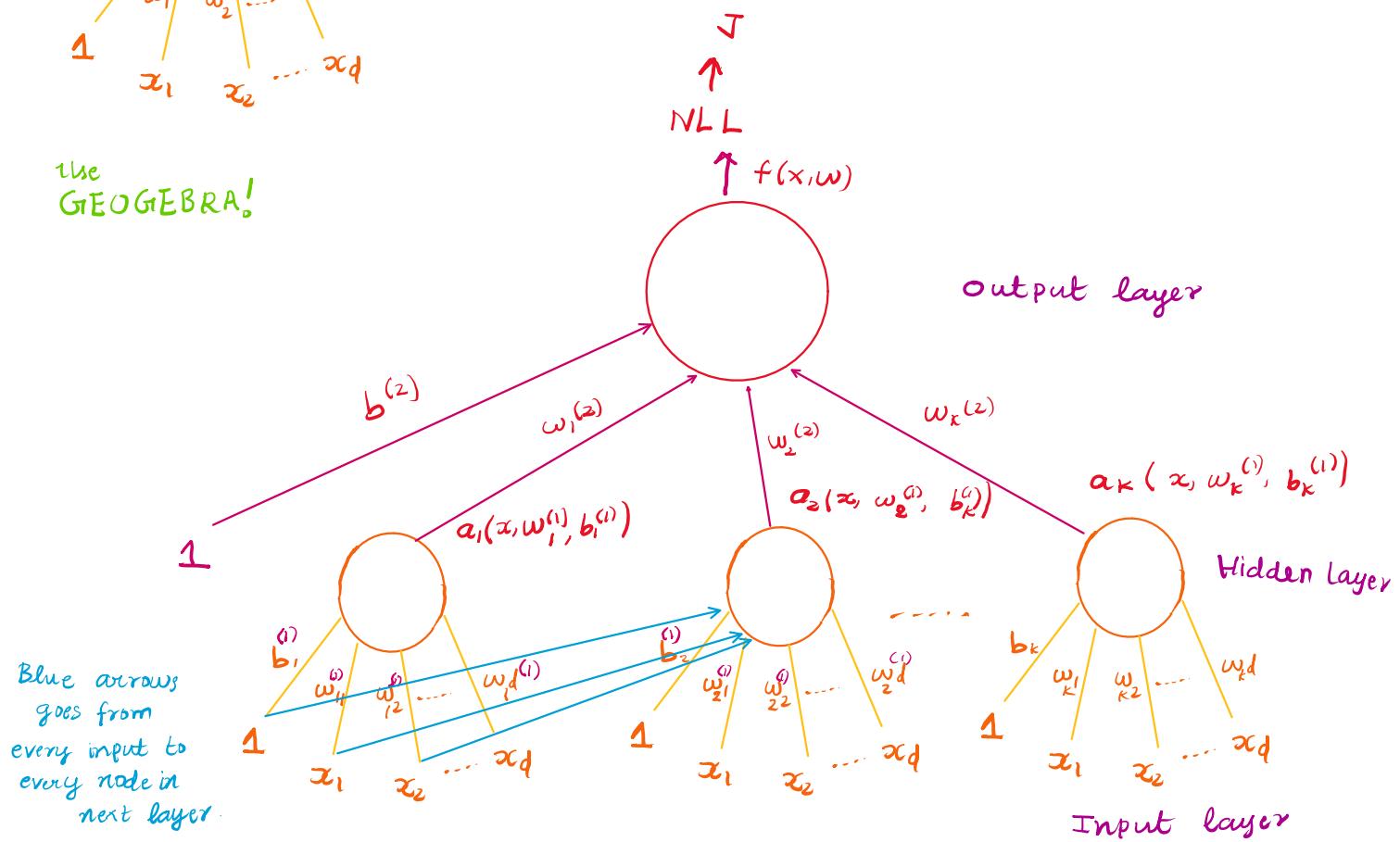
$$g = \sigma(\text{sigmoid})$$



$$NLL_i = -[y_i \log f(x_i, w) + (1 - y_i) \log(1 - f(x_i, w))]$$

(cross entropy)

the  
GEOGEBRA!

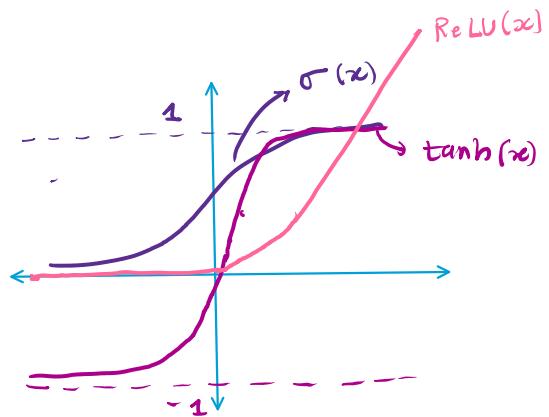


GOAL: To learn the  $w$ 's and  $b$ 's of every layer to minimize the loss.

### Popular Activation fns

i) Sigmoid  $\sigma(x) = \frac{1}{1 + e^{-x}}$

ii) Hyperbolic tan :  $\tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$   
(scaled sigmoid)



ii) Rectified linear units:

$$\text{ReLU}(x) = \max\{0, x\}$$

**Vanishing Gradient problems:** For deeper neural networks, the gradient vanishes for large 'x'.

Gradient doesn't vanish for ReLU fn, hence it's desirable for large x.

### Training NN:

Given a dataset  $D = \{(x_i, y_i)\}_{i=1}^n$  find the values of w's and b's that minimize the loss fn. (Fit forward NN)

STEP 1: Define a loss function  $J(w, b)$ . (eg: cross entropy loss)

STEP 2:  $J(\theta) = \sum_{i=1}^n l(NN(x_i, \theta), y_i)$  (here,  $f(x_i, \theta) = NN(x_i, \theta)$ )  
{ eg:  $l(NN(x_i, \theta), y_i) = [-y_i \log(NN(x_i, \theta)) + (1-y_i) \log(1-NN(x_i, \theta))]$  }

Once the loss fn is defined, we use SGD to optimize it.

### NN training algorithm:

- Inputs:  $NN(x, \theta)$ , training examples  $x_1, x_2, \dots, x_n$ , labels  $y_1, \dots, y_n$  and a loss function  $l$
- Randomly initialize  $\theta$
- Do until stopping criteria:
  - Pick randomly an example  $(x_i, y_i)$
  - Compute gradient of  $l$ ,  $\nabla_{\theta} l(x_i, y_i)$
  - $\theta_{t+1} \leftarrow \theta_t - \eta \nabla_{\theta} l$
- Return  $\theta_{t+1}$

Mini-batch GD, instead of SGD

choose a batch B  
of examples  
 $B \subseteq \{1, \dots, n\}$

$$\sum_{i \in B} \nabla_{\theta} l(x_i, y_i)$$

How will we compute  $\nabla_{\theta} l$  efficiently? Using backpropagation, which uses the chain rule of differentiation in a clever way!

Scalars:  $\frac{dg}{du} = \frac{du}{dx} \frac{dg}{dx}$ , where  $\bar{g} = g(u)$  and  $u = u(x)$

Scalars:  $\frac{dg}{dx} = \frac{du}{dx} \frac{dg}{du}$ , where  $\bar{g} = g(u)$  and  $u = u(x)$   
 ( $g, x, u$  are scalars)

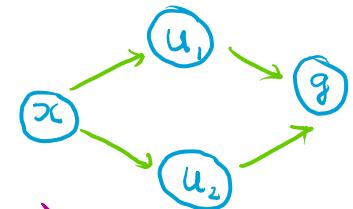


Vectors:  $\frac{\partial g}{\partial x} = \frac{\partial u}{\partial x} \cdot \frac{\partial g}{\partial u} + \frac{\partial u_2}{\partial x} + \frac{\partial g}{\partial u_2}$

( $g, x$  are scalars, but  $u = (u_1, u_2)$  is vector)

$$= \frac{\partial u}{\partial x} \cdot \frac{\partial g}{\partial u}$$

$$\left( \frac{\partial u_1}{\partial x}, \frac{\partial u_2}{\partial x} \right) \quad \left( \frac{\partial g}{\partial u_1}, \frac{\partial g}{\partial u_2} \right)$$



Assume  $g, x, u$  to be all vectors:

$$\frac{\partial g}{\partial x} = \frac{\partial u}{\partial x} \cdot \frac{\partial g}{\partial u}$$

where,

$$\frac{\partial u}{\partial x} = \begin{bmatrix} \frac{\partial u_1}{\partial x_1} & \frac{\partial u_2}{\partial x_1} & \dots & \frac{\partial u_k}{\partial x_1} \\ \frac{\partial u_1}{\partial x_2} & \dots & \dots & \frac{\partial u_k}{\partial x_2} \\ \vdots & \ddots & \ddots & \vdots \end{bmatrix}$$

and,

$$\frac{\partial g}{\partial u} = \begin{bmatrix} \frac{\partial g_1}{\partial u_1} & \dots & \frac{\partial g_p}{\partial u_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_1}{\partial u_k} & \dots & \frac{\partial g_p}{\partial u_k} \end{bmatrix}$$

Back propagation :  $\frac{\partial g_i}{\partial x_j} = \sum_{j=1}^k \frac{\partial u_j}{\partial x_j} \cdot \frac{\partial g_i}{\partial u_j}$