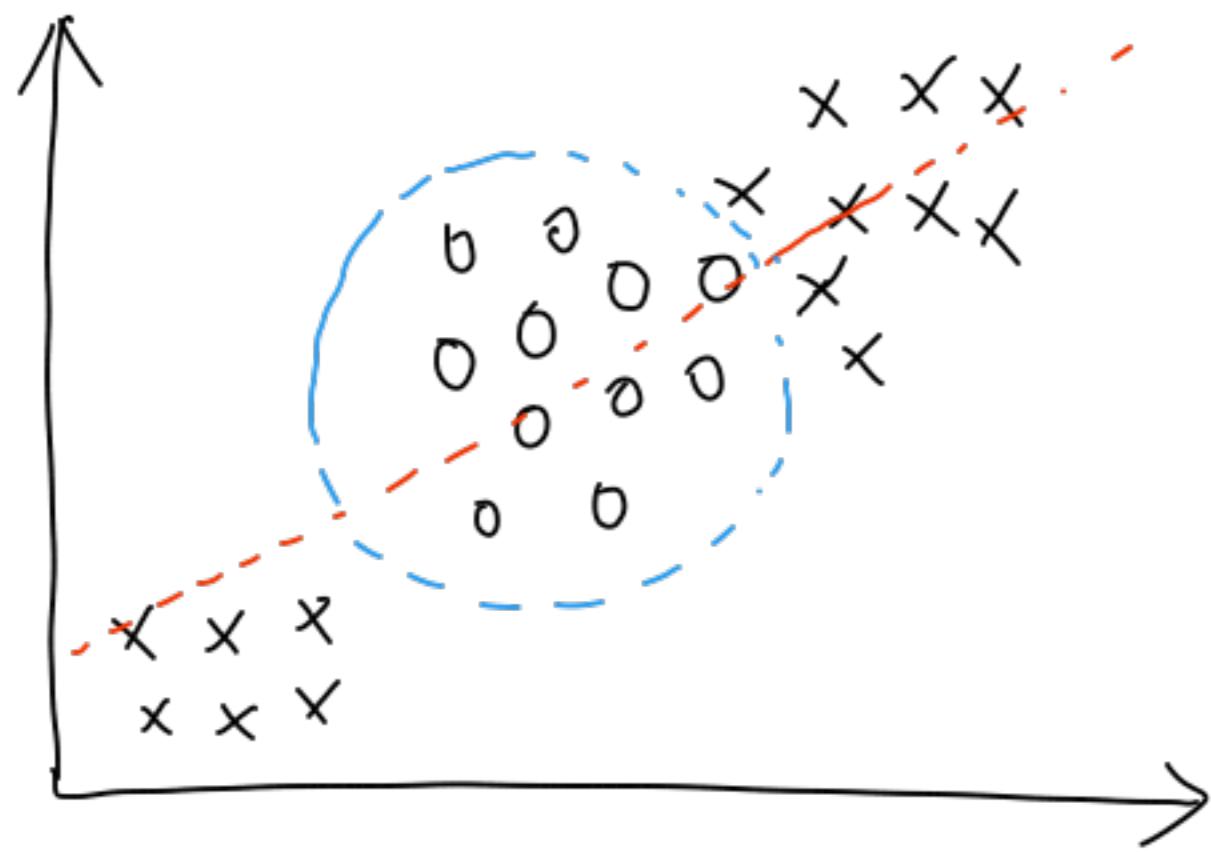# Lec 11: Neural Networks

Ex. 1



Vanilla LR

$$f(x, w) = \frac{1}{1 + e^{-w^T x}}$$

$$f_B(x, w) = \frac{1}{1 + e^{-w^T \phi(x)}}$$

$$\phi(x) = \begin{pmatrix} 1 & x_1 & x_2 & x_1^2 & x_2^2 \end{pmatrix}^T$$
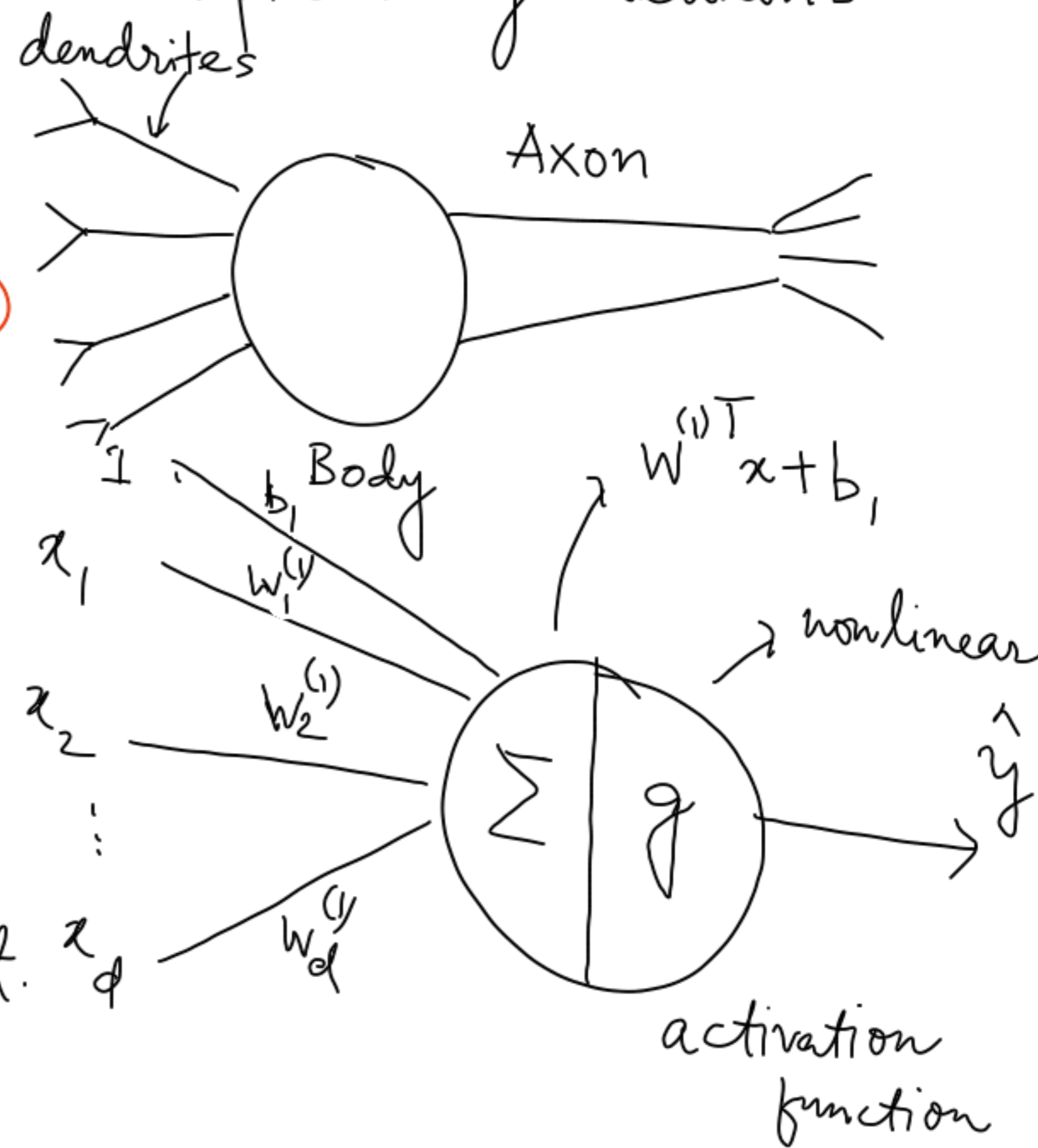
Philosophy of NN: Come up with some $\phi(x)$ without explicitly programming it.

Origin   1943 → 1960 →

1980's

Inspired by Neurons

dendrites
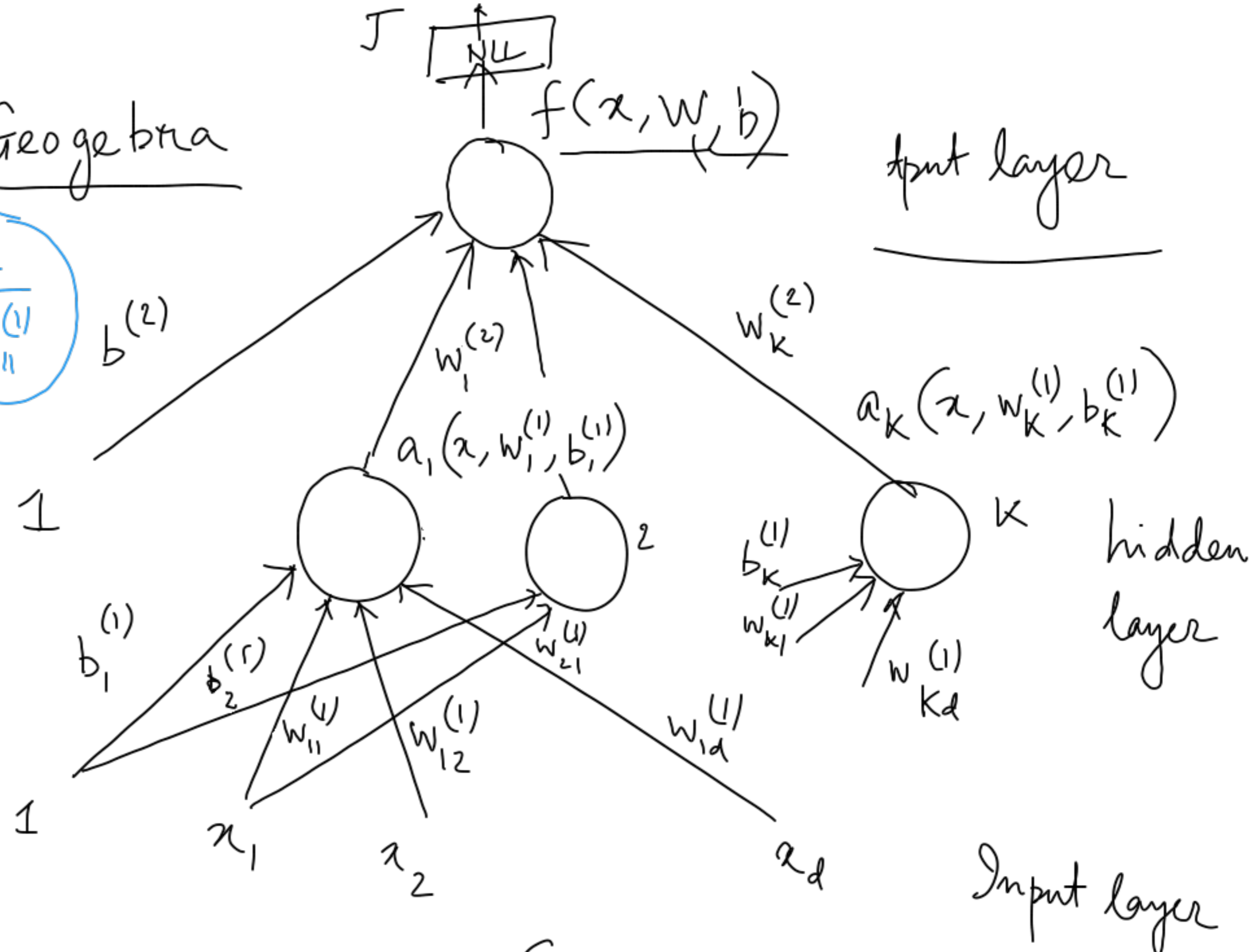
Axon

Body

$w^{(1)T} x + b_1$

$x_1$  1   $b_1$  $w_1^{(1)}$

$x_2$  $w_2^{(1)}$

$x_d$  $w_d^{(1)}$

$\sum$  $g$

nonlinear   $\hat{y}$

activation function

$\nabla_\theta \ell$

$\uparrow J$

NLL

$f(x,w)$

1

$b$   $w_1$   $w_2$   $w_d$

$x_1$   $x_2$   $x_d$

$$f(x,w) = g\left(W^T x + b\right)$$

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

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

cross-entropy

Geogebra

NN representation of Binary LR

$\dfrac{\partial \ell}{\partial w_{11}^{(1)}}$   $b^{(2)}$

$J$   NLL

$f(x, W, b)$

tput layer

$w_k^{(2)}$

$w_1^{(2)}$   $a_k(x, w_k^{(1)}, b_k^{(1)})$

$a_1(x, w_1^{(1)}, b_1^{(1)})$

1   2   $K$   hidden layer

$b_k^{(1)}$

$w_{k1}^{(1)}$

$b_1^{(1)}$

$b_2^{(1)}$   $w_{21}^{(1)}$

$w_{11}^{(1)}$   $w_{12}^{(1)}$   $w_{1d}^{(1)}$   $w_{kd}^{(1)}$

1   $x_1$   $x_2$   $x_d$   Input layer

Goal: to learn the W's and b's of every layer to minimize the loss.

# Popular activation functions
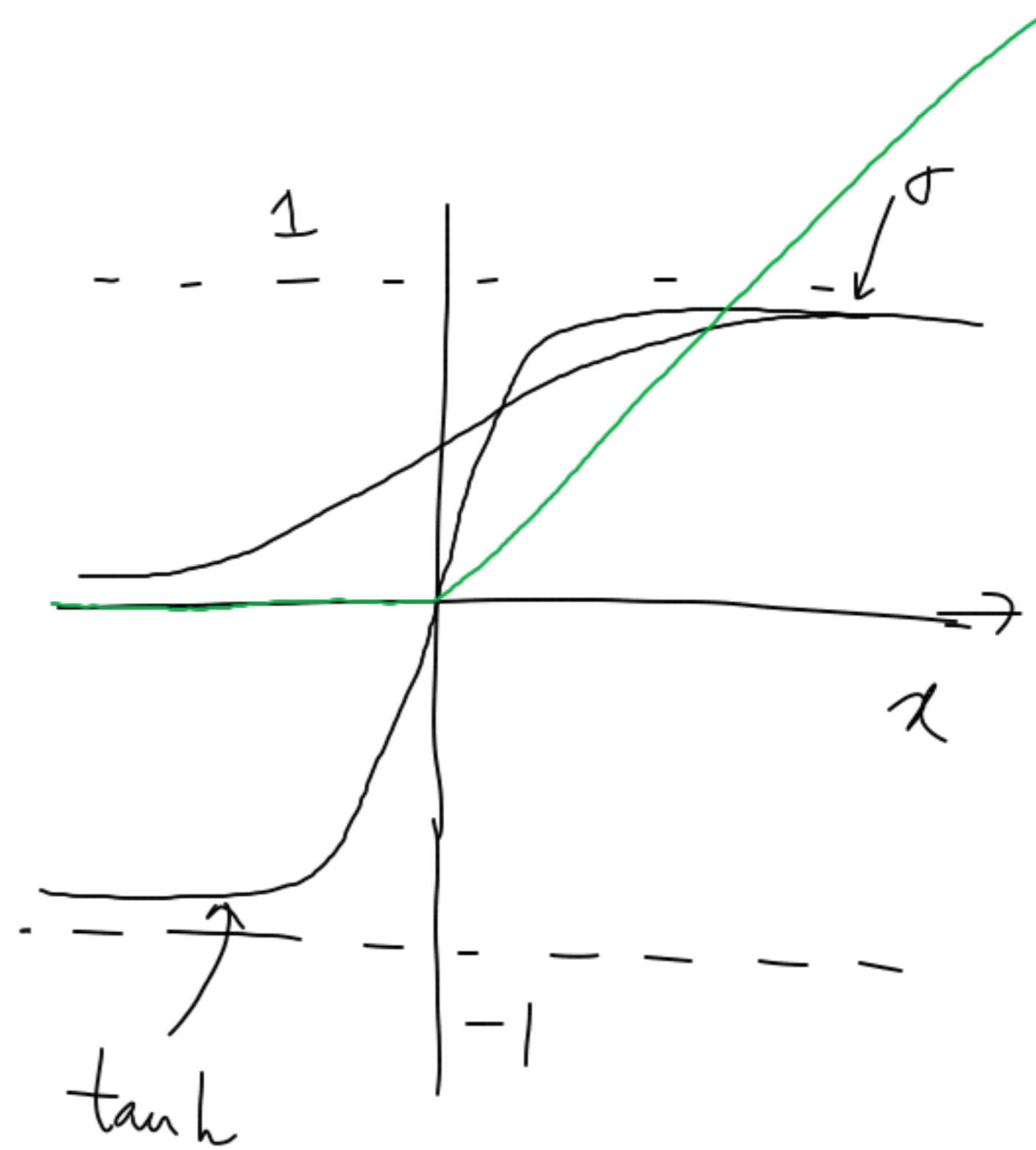
Sigmoid: $\sigma(x) = \dfrac{1}{1 + e^{-x}}$

Hyperbolic tan: $\tanh(x) = \dfrac{e^{2x} - 1}{e^{2x} + 1}$

$\underset{\text{(Scaled sigmoid)}}{} = 2\sigma(2x) - 1$

Rectified Linear Units:

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

gradient does not vanish for ReLU.



tanh

## Vanishing gradient problem:

for deeper NNs, the gradient vanishes for large $x$.

$D = \left( x_i, y_i \right)_{i=1}^{m} \rightarrow$ given this dataset

find the values of $W$'s and $b$'s that

minimize the loss function $\boxed{\text{Training NN}}$

Feed forward NN

e.g. $\qquad \nearrow P(y_i = 1 | x_i, \theta)$

$f(x_i, \theta) = NN(x_i, \theta)$

$\qquad \searrow \text{Soft max}(x, \theta)$

For a full blown NN

Step 1: Define a loss function $J\left(\overset{\theta}{\overbrace{W, b}}\right)$, $\left[\text{e.g. cross entropy loss}\right]$

Step 2:

$$J(\theta) = \sum_{i=1}^{n} \overset{\searrow \text{ell}}{\ell}\left( NN(x_i, \theta), y_i \right)$$

$$\ell\left( NN(x_i, \theta), y_i \right) = -\left[ y_i \log\left( NN(x_i, \theta) \right) + (1 - y_i) \log\left( 1 - NN(x_i, \theta) \right) \right]$$

example.

## NN training algorithm

- Inputs: $NN(x, \theta)$ , training examples $x_1, \dots, x_n$, labels $y_1 \dots, y_n$

  and a loss function $\ell$

- Randomly initialize $\theta$     mini-batch $B \subseteq \{1, \dots, n\}$

- do until stopping criteria:

              a batch $B$ of examples.

      Pick randomly an example $(x_i, y_i)$

      Compute gradient of $\ell$ , $\nabla_\theta \ell(x_i, y_i)$

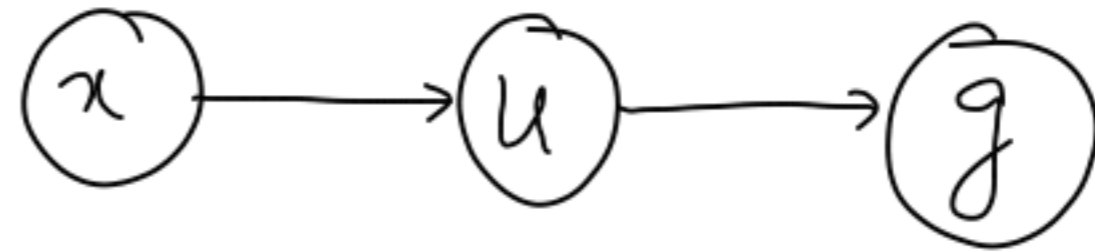      $\theta_{t+1} \leftarrow \theta_t - \eta \nabla_\theta \ell$

- Return $\theta_{t+1}$

How to compute $\nabla_\theta \ell$ efficiently?

via Backpropagation. → uses the chain rule of differentiation in a clever way.

Scalars:
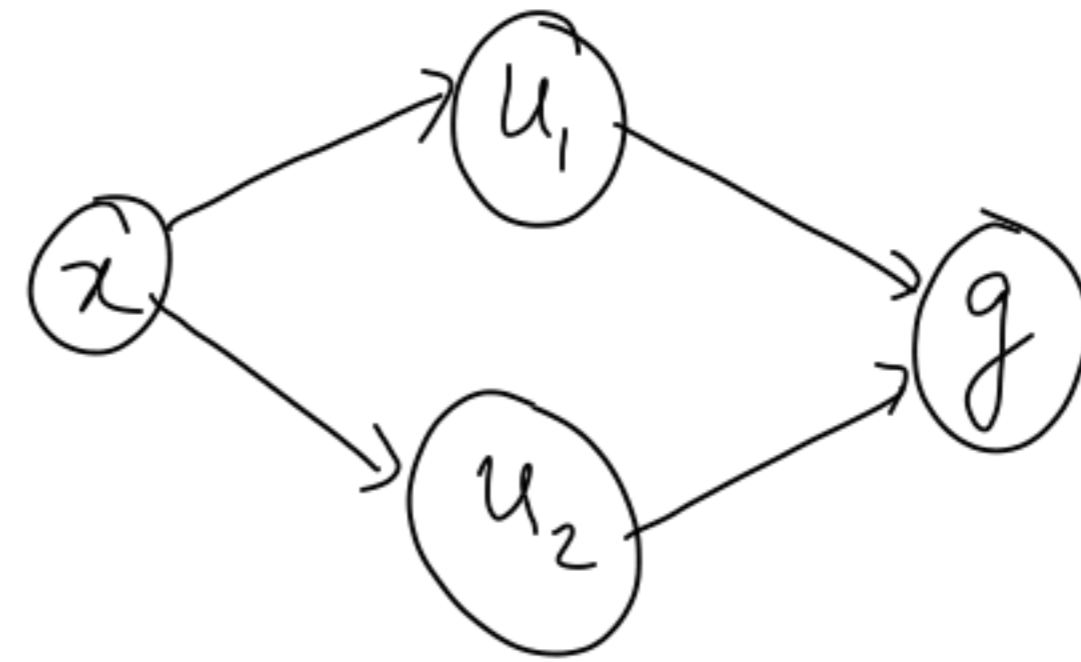$$\frac{dg}{dx} = \frac{du}{dx} \cdot \frac{dg}{du}$$
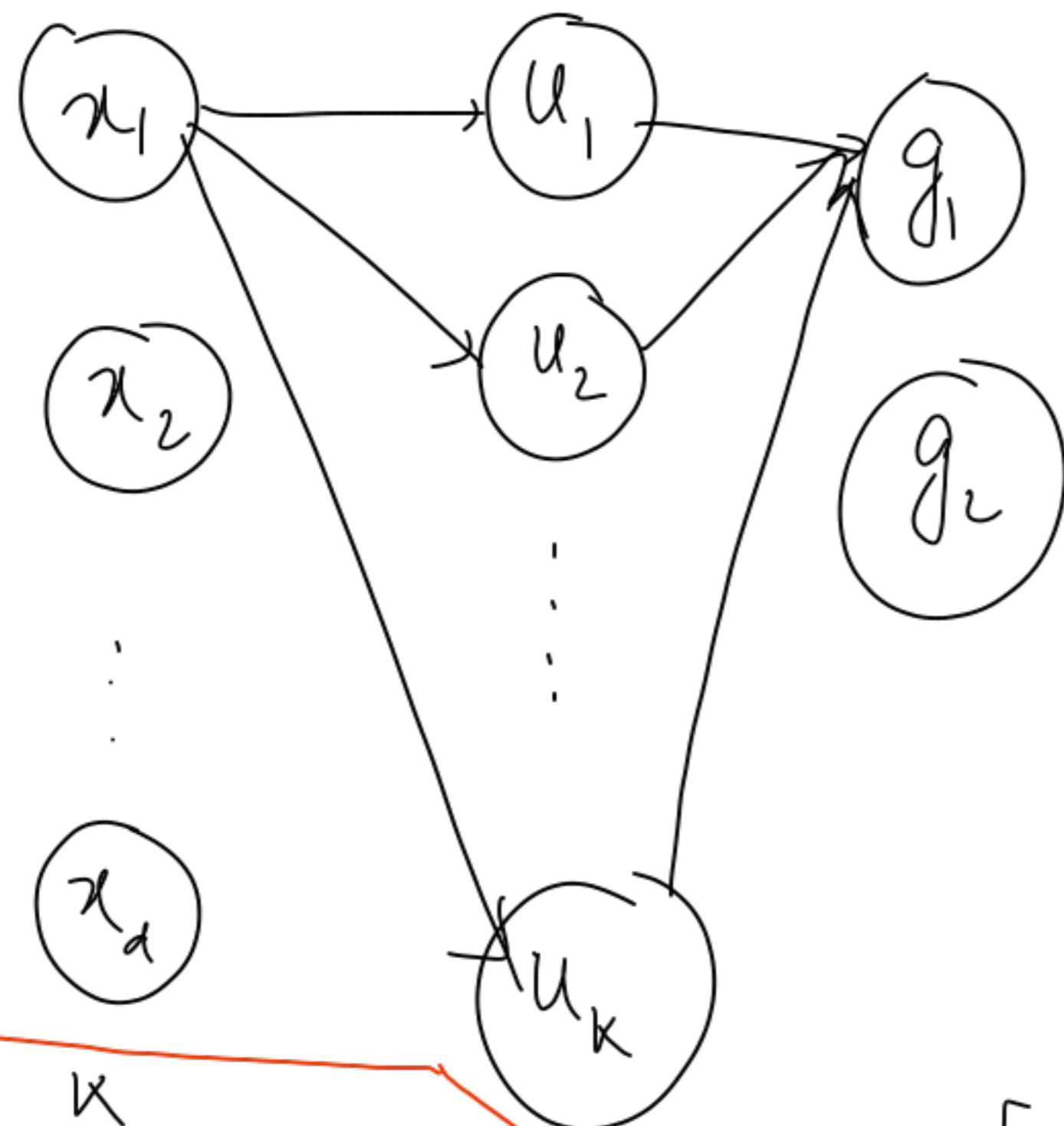
Vectors:
$$\frac{\partial g}{\partial x} = \frac{\partial u_1}{\partial x} \cdot \frac{\partial g}{\partial u_1} + \frac{\partial u_2}{\partial x} \cdot \frac{\partial g}{\partial u_2}$$

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

$$\begin{bmatrix} \frac{\partial u_1}{\partial x} & \frac{\partial u_2}{\partial x} \end{bmatrix}$$

$$\frac{\partial g}{\partial \underline{u}} = \begin{bmatrix} \frac{\partial g}{\partial u_1} & \frac{\partial g}{\partial u_2} \end{bmatrix}$$

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

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

$$\boxed{\frac{\partial g_1}{\partial x_1} = \sum_{j=1}^{K} \frac{\partial u_j}{\partial x_1} \cdot \frac{\partial g_1}{\partial u_j}}$$

Backpropagation.

$$\frac{\partial g}{\partial u} = \begin{bmatrix} \frac{\partial g_1}{\partial u_1} & \cdots & \frac{\partial g_m}{\partial u_1} \\ \frac{\partial g_1}{\partial u_2} & & \vdots \\ \vdots & & \frac{\partial g_m}{\partial u_k} \\ \frac{\partial g_1}{\partial u_k} & & \end{bmatrix}$$