Recall that a basic approach for constructing non-linear function classes in using basis functions,

$$f(x, a) = \sum_{\ell=1}^{m} a_\ell \phi_\ell(x) = a^T \phi(x).$$

where $\phi(x) = [\phi_1(x), \cdots, \phi_m(x)]^T$ is a set of basis functions that is believed to capture important information regarding the input, and $a = [a_1, \cdots, a_m]^T$ is a coefficient vector applied on $\phi(x)$.

for example, in polynomial regression, we assume $\phi_\ell(x) = x^{\ell-1}$ and hence $f(x, a) = \sum_{\ell=0}^{m-1} a_\ell x^\ell$.

## Neural Networks

Approach to non-linear Regression:

### Fixed basis function:

$$f(x; w) = \sum_i w_i \phi_i(x)$$

e.g. polynomial regression: $f(x, w) = w_0 + w_1 x + w_2 x^2 + \ldots$

### Adaptive basis functions:

* Kernel method: $f(x, w) = \sum_{i \in data} w_i k(x, x^{(i)})$

* Neural Network: $f(x; w, v) = \sum_i w_i \phi(x; v_i)$

$$\underset{i^{th}\ neuron}{\phi_i(x)} = \underset{\substack{\text{coefficient of the basis functions.}\\ \text{(weight of neuron } i)}}{\phi(x, v_i)}$$

**Objective:** $\underset{w, v}{\text{Min}}\ E_D\left[(y - f(x; w, v))^2\right]$.

**Basis function (neuron) form:** Each neuron $\phi_i$ is assumed to have the following form:

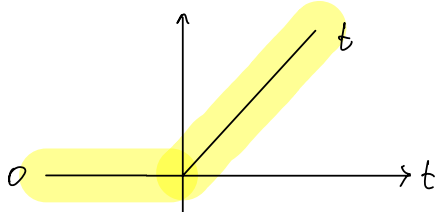$$\underset{\substack{\text{basis}\\ \text{function}\\ \text{aka}\\ \text{neuron}}}{\phi_i(x)} = \underset{\substack{\text{activation}\\ \text{function}\\ \text{(non-linear).}}}{\sigma}\left(\sum_{\ell=1}^{d} v_\ell x_\ell + v_0\right), \quad \text{where } x = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(d)} \end{bmatrix}, \quad v = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix}$$

each neuron $\phi_i$ has an activation function, each which is composed of summation of $v_\ell x_\ell + v_0$, $\ell = 1$ to $d$.
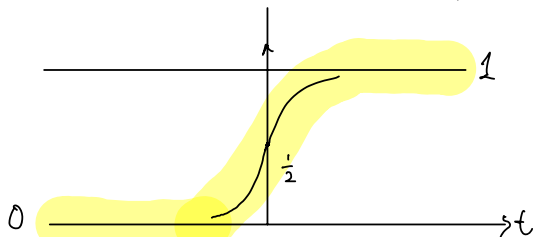
Examples of activation function:

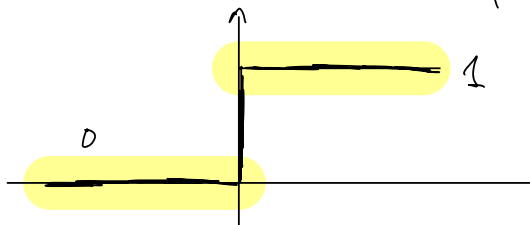1) ReLU (Rectified Linear Unit): $\sigma(t) = \max(0, t)$.

$$\sigma(t) = \begin{cases} 0, & t \geq 0 \\ t, & t > 0. \end{cases}$$

2) Sigmoid: $\sigma(t) = \dfrac{\exp(t)}{1 + \exp(t)} = \dfrac{1}{1 + \exp(-t)}$

3) Indicator:

$$\sigma(t) = \mathbb{I}(x > 0) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

neuron $\in \{1, \dots m\}$

$d$ = data point

linear coefficient $\in \{1, \dots m\}$

weight $\ell$ of $d$ for neuron $i$:
$v_{i\ell} \in \{0, \dots d\}$

$m \sim$ complexity of NN.

$$\underbrace{f(x; w, v)}_{\text{neural network}} = \underbrace{\sum_{i=1}^{M} w_i}_{\substack{\text{summing all neurons} \\ \text{from 1 to } m.}} \sigma \Big( \underbrace{\sum_{\ell=1}^{d} v_{i\ell} x_\ell + v_{i0}}_{\substack{\text{each neuron has its own components,} \\ \text{from 1 to } d.}} \Big),$$

where $v_i = \begin{bmatrix} v_{i0}, & v_{i1}, & \dots, & v_{id} \end{bmatrix}^T$ // vector

$v = \begin{bmatrix} v_1, & \dots, & v_m \end{bmatrix}$ // matrix

Loss function : $\min\limits_{w, v} E_D\left[(y - f(x; w, v))^2\right]$.

use GD or SGD.

Learning a neural network is a nonconvex optimization problem! No closed form, can only find local minima.

Example: A NN with 2 neurons:
$$f(x, v) = \sigma(x - v_1) + \sigma(x - v_2)$$
$$L(v) = E_D\left[(y - \sigma(x - v_1) - \sigma(x - v_2))^2\right].$$
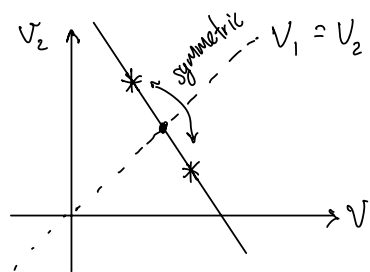
Assume $[v_1^*, v_2^*]$ is optimal.
Then $[v_2^*, v_1^*]$ is also optimal.

If $L(v)$ is convex,
then $[\overline{v}^*, \overline{v}^*]$ is also optimal,
where $\overline{v}^* =$ average of the 2.
$$= \frac{v_1^* + v_2^*}{2}$$



non convex

mostly

usually not

convex

## Gradient Descent: To solve the optimization.

$$\min_{\omega, v} L(\omega, v) : \begin{cases} \omega^{new} = \omega^{old} - \epsilon \nabla_\omega L(\omega, v) \\ v^{new} = v^{old} - \epsilon \nabla_v L(\omega, v) \end{cases}$$

- Initialization is critical. Because where you initialize decides which local min you end up with.

- Initialize $v$ to be "random small values".
- Initialize the neurons differently.
- Be careful with initialization of neurons and step sizes.
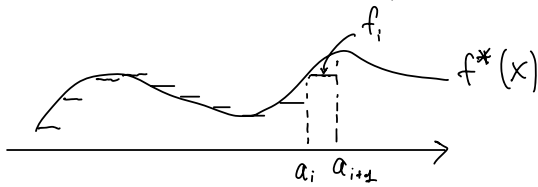
## Example: Digit Classification:

X = matrix                          x = vector



1 or 7?                             $d \times 1$

$\xrightarrow{reshape}$     $\xrightarrow[f]{function}$   $\{ 0, \cdots, 9 \}$



$$\sum_{i=1}^{m} \omega_i \sigma \left( v_i^T x - v_o \right)$$

$d \times 1$

This depends on the activation function:

weight $v_i$ is estimated by minimizing the loss.

Take inner product of X (image) with $v_i$ (weights) and see which one has the least loss. That is the closest template.

If $\sigma = \mathbb{I}$ :    $\sum_{i-1}^{m} w_i \, \mathbb{I}\left(v_i^T x \geq v_o\right)$    // if inner product $> v_o$,
                                                            then neuron fires.

## Theoretical Results

If you have infinite neurons, $f(x; w, v)$ can approximate
arbitrary nonlinear functions:

If $m$ is large enough, then the NN $f(x, w, v)$
can approximate any continous function.



can approximate using step (indicator) function:

$$\sum_{i=1}^{m} f_i \, \mathbb{I}\left(a_i \leq x \leq a_{i+1}\right)$$

$$= \sum_{i=1}^{m} f_i \left(\mathbb{I}(x \geq a_i) - \mathbb{I}\left(x \geq a_{i+1}\right)\right)$$

$$= \sum_{i=1}^{m} \left(f_i - f_{i-1}\right) \mathbb{I}\left(x \geq a_i\right)$$

$\Rightarrow$ There must exist some parameter that can approximate
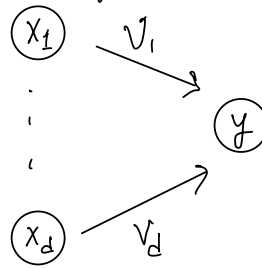   a continous function arbitrarily well.

using graph representation is easier.

# Graphical Representation of NN:

Similar to Gaussian Markov Field.

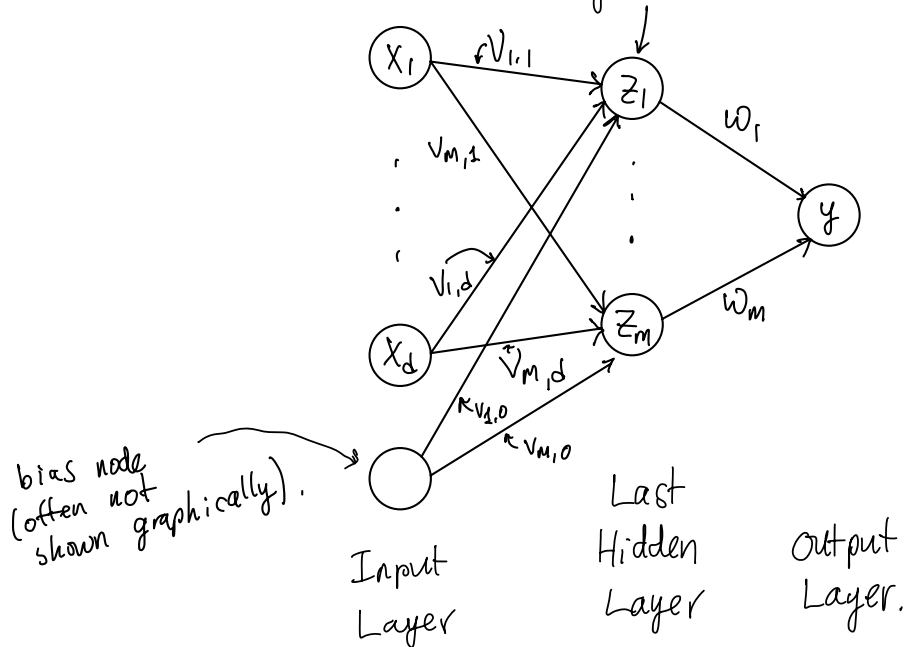We can illustrate this using linear regression:

$$y = \sum_{l=1}^{d} v_l x_l$$

← data points

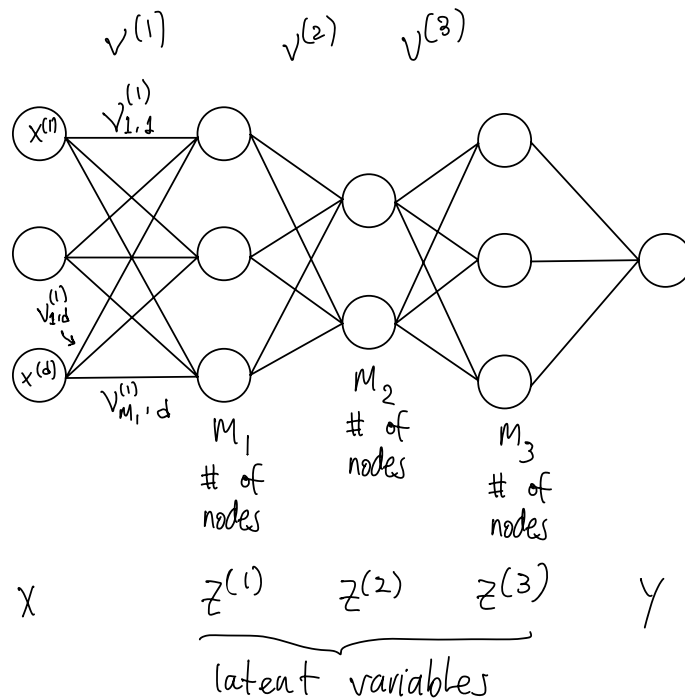$$y = \sum_{i=1}^{M} w_i \, \sigma\left( \underbrace{\sum_{l=1}^{d} v_{i,l} x_l + v_0}_{\triangleq z_i \; = \; \text{output of each neuron.}} \right)$$

output of previous layer's neuron.

bias node (often not shown graphically).

Input Layer

Last Hidden Layer

Output Layer.

If multiple hidden layers : Deep neural network.

$V^{(1)}$    $V^{(2)}$    $V^{(3)}$



$x$    $\underbrace{z^{(1)} \quad z^{(2)} \quad z^{(3)}}_{\text{latent variables}}$    $y$

$$z^{(1)} = \Phi\left(x, V^{(1)}\right) = \begin{bmatrix} \sigma\left(\left(V_1^{(1)}\right)^T x\right) \\ \vdots \\ \sigma\left(\left(V_{M_1}^{(1)}\right)^T x\right) \end{bmatrix}$$

⟵ first neuron in $z^{(1)}$.

size: $M_1 \times d$.

⟵ last neuron in $z^{(1)}$.

$$z^{(2)} = \Phi\left(z^{(1)}, V^{(2)}\right) = \begin{bmatrix} \sigma\left(\left(V_1^{(2)}\right)^T z^{(1)}\right) \\ \vdots \\ \sigma\left(V_{M_2}^{(2)}\right)^T z^{(1)} \end{bmatrix}$$

⟵ first neuron in $z^{(2)}$.

size: $M_2 \times M_1$.

⟵ last neuron in $z^{(2)}$.

$$z^{(3)} = I\left(z^{(2)}, V^{(3)}\right) = \cdots$$

$$y = \sum_{i=1}^{M_3} W_i \, z_i^{(3)}$$

Backpropagation = <mark>calculation of gradients of DNNs</mark>.

Backpropagating the error from output to inputs.

Consider an example DNN where each layer has 1 neuron:

$$x \xrightarrow{v^{(1)}} z^{(1)} \xrightarrow{v^{(2)}} z^{(2)} \cdots \xrightarrow{v^{(m)}} y$$

$$L(v) = E_D\left[\left(y - f(x, v^{(1)} \cdots v^{(m)})\right)^2\right]$$

$$\nabla_{v^{(1)}} L(v) = 2 E_D\left[\left(f(x,v) - y\right) \frac{df(x,v)}{dv^{(1)}}\right] \quad \text{// chain rule.}$$

where

$$\frac{df(x,v)}{dv^{(1)}} = \frac{dy}{dv^{(1)}} = \left(\frac{dy}{dz^m}\right)\left(\frac{dz^m}{dz^{m-1}}\right) \cdots \left(\frac{dz^3}{dz^2}\right)\left(\frac{dz^2}{dz^1}\right)\left(\frac{dz^1}{dv^1}\right) \quad \text{// chain rule.}$$