# Some Interesting Interpretations (MLP)

(*) These are **Feedforward Networks**. If we include feedback, these become **Recurrent Networks**.

(*) To remove the restriction of linear models, we consider restricted non-linear models

One possible (Kernels): instead of $\underline{x}$, consider $\phi(\underline{x})$
approach     $\phi(\cdot) \rightarrow$ possible non-linear mapping

Three possible methods

(1) General RBF kernel: SVM-based 'black box' methods use a generic RBF kernel

(2) Manually engineer $\phi(\underline{x})$: difficult! does not generalise across domains e.g., speech & computer vision domains

(3) Deep learning: typically learn $\phi(\underline{x})$

$$\underline{y}(\,) = \underline{f}(\underline{x}) = \underline{f}(\underline{x}, \underline{w})$$

Earlier linear approach                    linear in $\phi(\cdot)$
$\underline{w}^T \underline{x} \longrightarrow \widetilde{\underline{w}}^T \phi(\underline{x})$
                                           itself is possibly non-linear
(homogeneous representation)

$$\phi(\underline{x}) = \phi(\underline{x}, \underline{\theta}) \leftarrow \text{parameters}$$

We can learn $\phi(\cdot)$ from a broad class of functions
This generalises the 1st & 2nd methods. How?
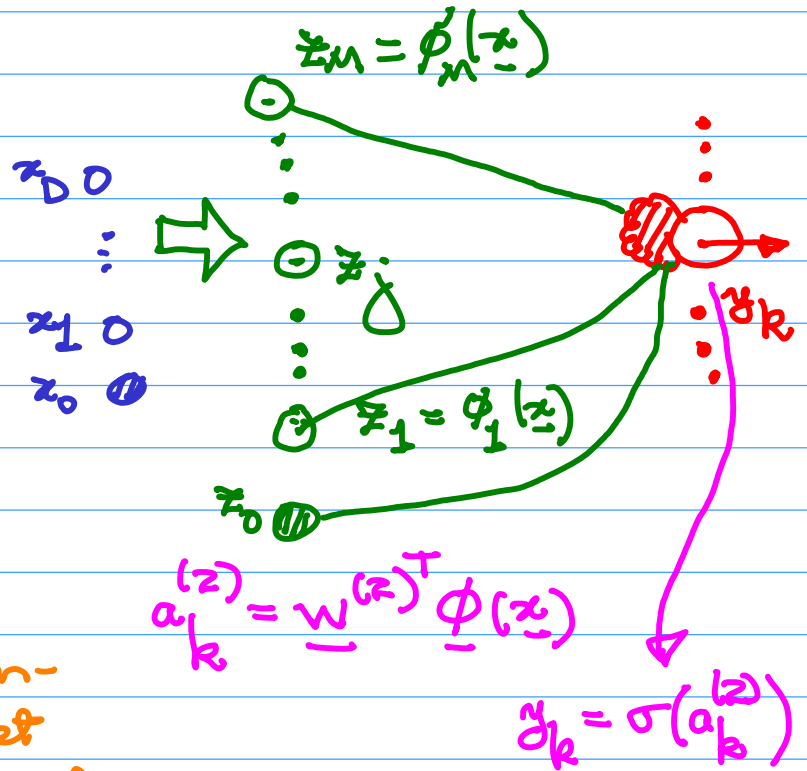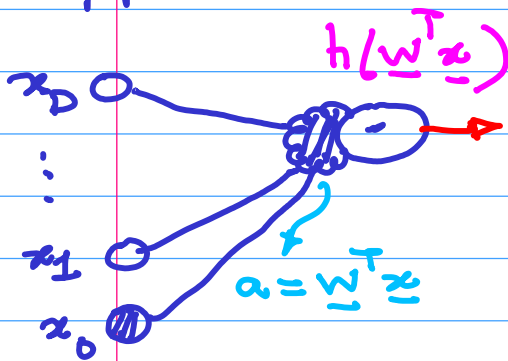   1st $\longrightarrow$ generic, very broad family
   2nd $\longrightarrow$ humanknowledge not to find the exact right function $\longrightarrow$

but to find the right family of functions for $\underline{\phi}(\underline{x})$

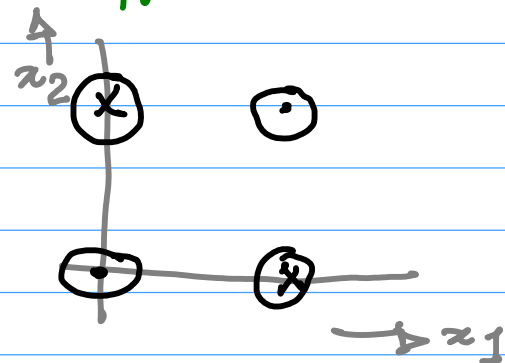**Main Idea** "$\underline{\phi}(\underline{x})$ defines a hidden layer"

Instead of the earlier approach

$h(\underline{w}^T \underline{x})$

$x_D$ ◯
⋮
$x_1$ ◯
$x_0$ ◉

$a = \underline{w}^T \underline{x}$

$z_M = \phi_M(\underline{x})$

$x_D$ ◯
⋮
$x_1$ ◯
$x_0$ ◉

⟹

◯ $z_j$

◯ $z_1 = \phi_1(\underline{x})$

$z_0$ ◉

$y_k$

$a_k^{(2)} = \underline{w}^{(2)T} \underline{\phi}(\underline{x})$
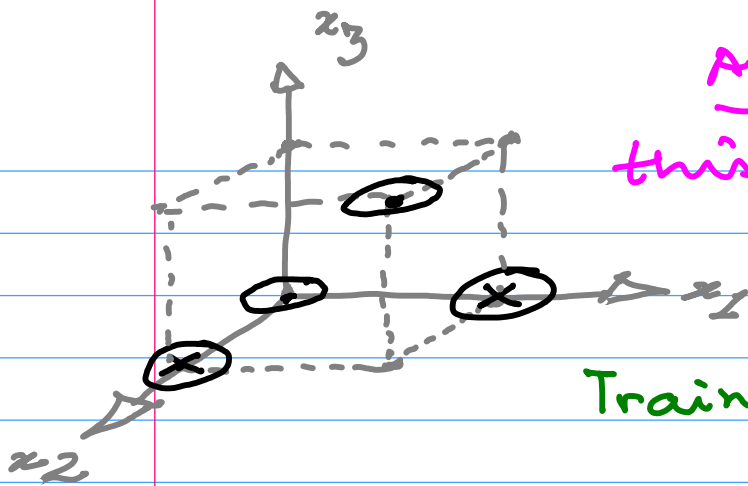
$y_k = \sigma(a_k^{(2)})$

We transform the inputs (possibly in a non-linear manner) to get the hidden layers, from where we take a linear combination of $\underline{\phi}(\underline{x})$ in the output layer

**the X-OR problem (again!)** what was our previous approach to solve it? Handcrafted kernel / feature transformation approach

| $x_2$ | $x_1$ | new feature $x_3 = x_2 x_1$ | $y = x_2 \oplus x_1$ |
|-------|-------|------|------|
| 0 | 0 | 0 | 0 → ◦ |
| 0 | 1 | 0 | 1 → ✗ |
| 1 | 0 | 0 | 1 → ✗ |
| 1 | 1 | 1 | 0 → ◦ |

$x_2$

✗    ⊙

⊙    ✗

→ $x_1$

**Approach here!** To formulate this as a regression problem and use mean square error loss.

Training set $\underline{x} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}$

$$X = \begin{bmatrix} x_{(1)} & x_{(2)} & x_{(3)} & x_{(4)} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$\left(\underline{t}\right) = \begin{bmatrix} t_{(1)} & t_{(2)} & t_{(3)} & t_{(4)} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}$$

MSE loss $J(\underline{w}) \triangleq \frac{1}{4} \sum_{n=1}^{4} \left[ y(x_{(n)}, \underline{w}) - t_{(n)} \right]^2$

$\hookrightarrow$ parameters

linear model

$$y(\underline{x}) = y(\underline{x}, \underline{w}) = \underline{w}^T \underline{x} = \underline{w}^T \underline{x} + b$$

$\underbrace{\underline{w}^T \underline{x}}_{\text{homogeneous}}$     $\underbrace{\underline{w}^T \underline{x} + b}_{\text{non-homogeneous}}$     or $w_0$