

# Data Mining: Artificial Neural Network

---

---

## Lecture Notes for Chapter 3

### Data Mining

<https://ml-graph.github.io/winter-2025/>

Yu Wang, Ph.D.

[yuwang@uoregon.edu](mailto:yuwang@uoregon.edu)

Assistant Professor

Computer Science

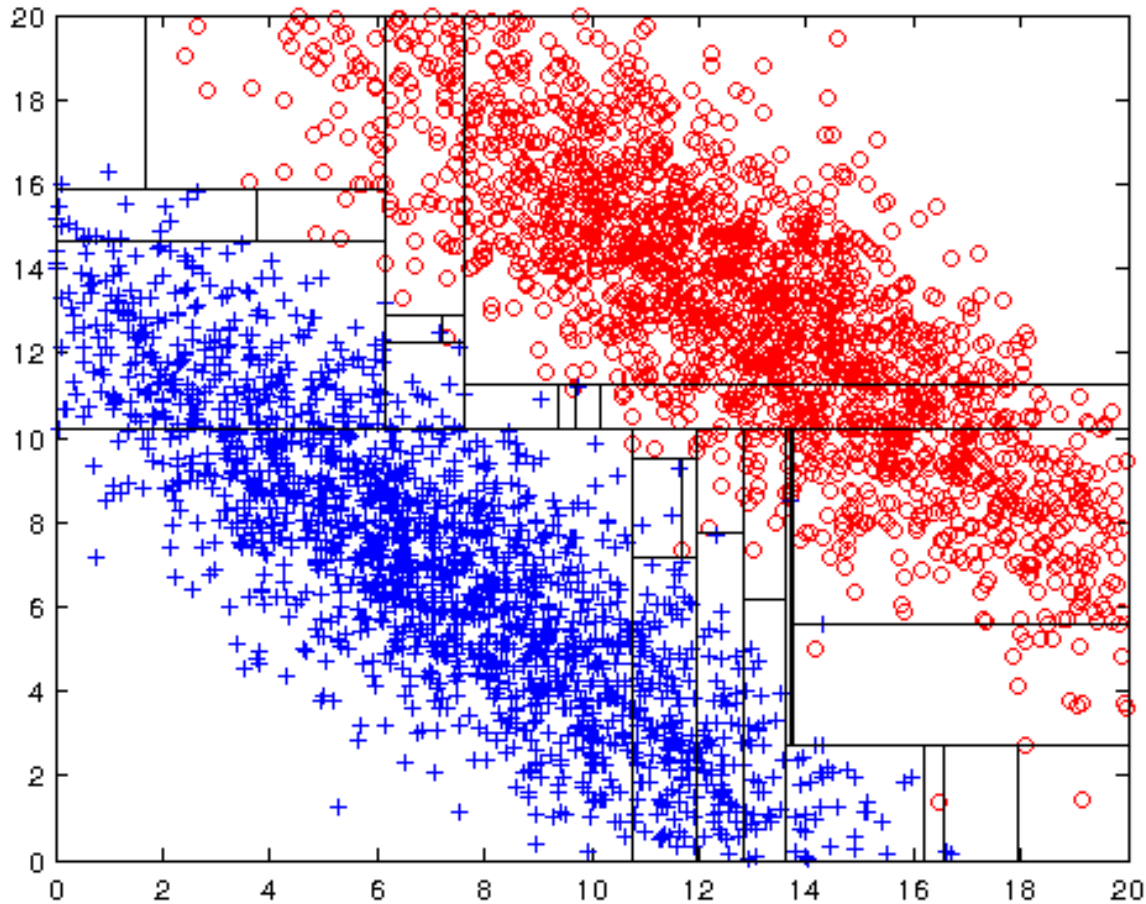
University of Oregon

CS 453/553 – Winter 2025

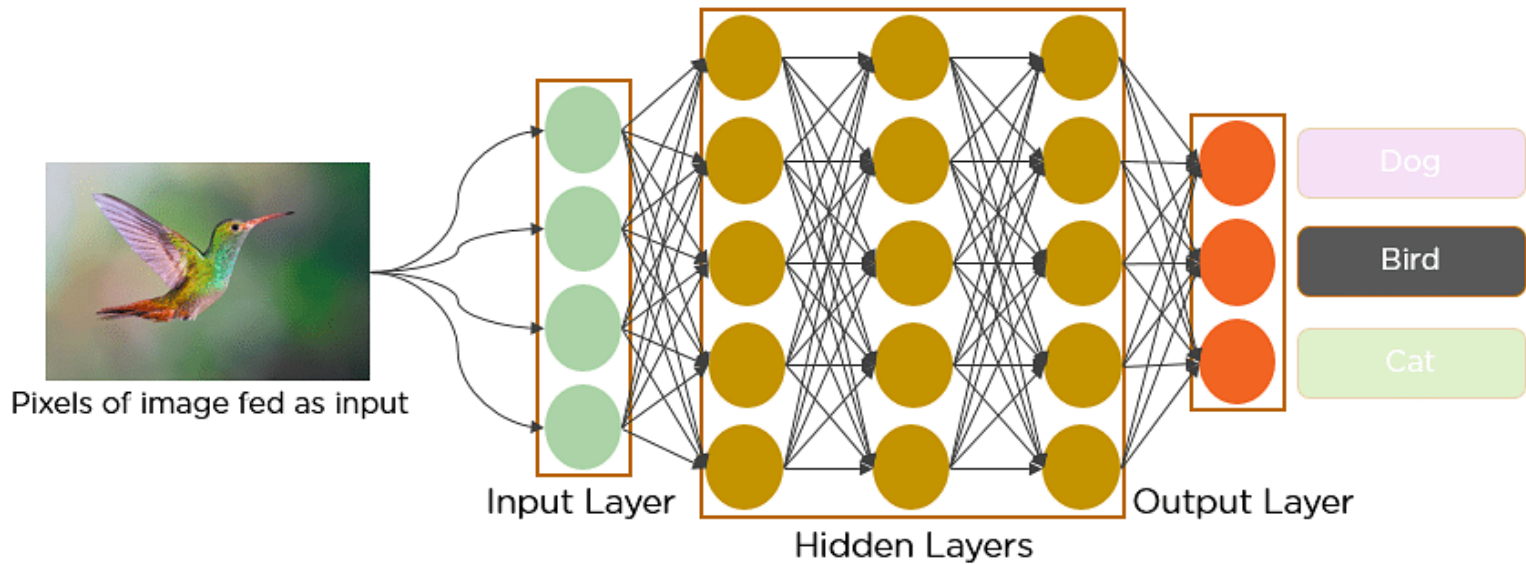
**Course Lecture is very heavily based on  
“Introduction to Data Mining”  
by Tan, Steinbach, Karpatne, Kumar**

# Limitations of Decision-Tree and many other model

---



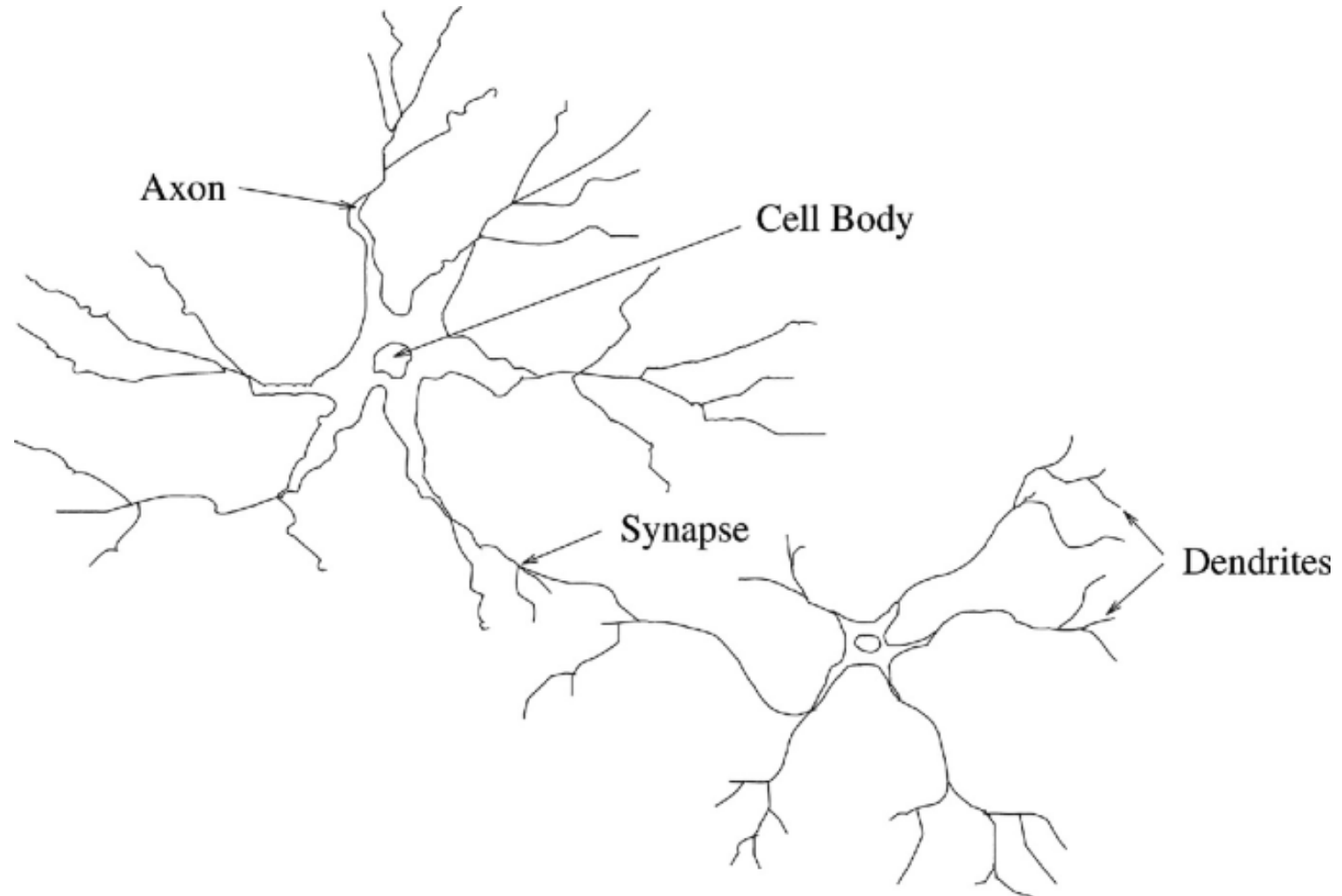
# Real-world Intuition



# Basic Architecture of Perceptron

---

---

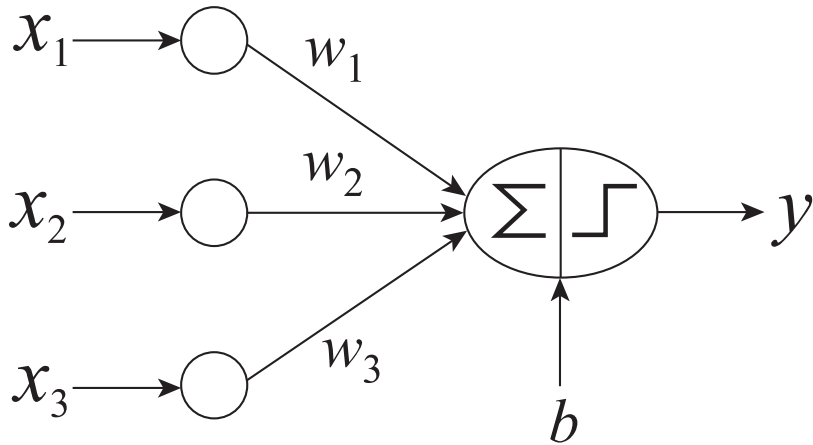


# Artificial Neural Networks (ANN)

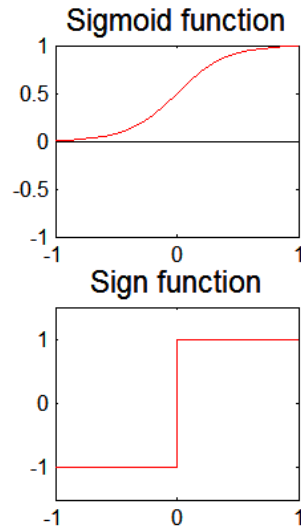
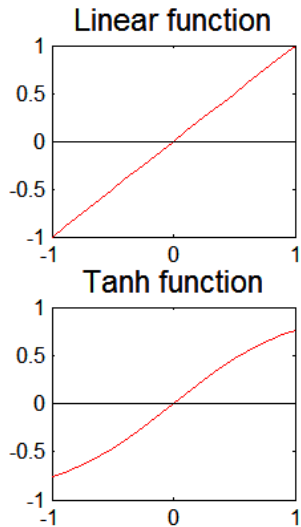
---

- **Basic Idea: A complex non-linear function can be learned as a composition of simple processing units**
- **ANN is a collection of simple processing units (nodes) that are connected by directed links (edges)**
  - Every node receives signals from incoming edges, performs computations, and transmits signals to outgoing edges
  - Analogous to *human brain* where nodes are neurons and signals are electrical impulses
  - Weight of an edge determines the strength of connection between the nodes
- **Simplest ANN: Perceptron (single neuron)**

# Basic Architecture of Perceptron



$$y = \sigma(w^T x + b)$$



**What happens if there is no nonlinear activation?**

# Linear Regression

---

- Data -  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$   $f(x) = xw + b$

- Regression – Find  $f$  that minimizes our uncertainty about  $y$  given  $x$

$$y = f(x) + n$$

- Minimizing Mean Squared Error = Minimizing Negative Log-Likelihood

$$\operatorname{argmin}_f \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}))^2$$

# Linear Regression

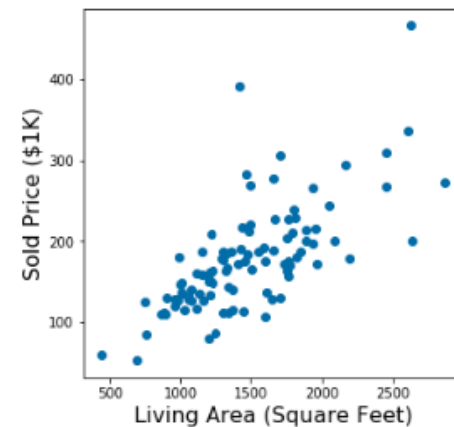
$$\operatorname{argmin}_w \frac{1}{N} \sum_{i=1}^N (y^{(i)} - w^T \hat{x}^{(i)})^2 = \operatorname{argmin}_w \frac{1}{N} \|y - Xw\|^2$$

Loss/Cost Function

Where

- $y = [y^{(1)}, \dots, y^{(N)}]^T \in \mathbb{R}^{N \times 1}$  and
- $X = [\hat{x}^{(1)}, \dots, \hat{x}^{(N)}]^T \in \mathbb{R}^{N \times (d+1)}$  (here  $d = 1$ )
- $w = [w_0, w_1, \dots, w_d]^T \in \mathbb{R}^{d+1}$

$$\operatorname{argmin}_f \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}))^2$$





# Linear Regression

---

---

$$J(w) = \frac{1}{N} \sum_{i=1}^N \left( y^{(i)} - w^T \hat{x}^{(i)} \right)^2$$

**Compute the minimum value? How to do it in Math?**

**Find points where gradient = 0**

# Linear Regression

---

$$J(w) = \frac{1}{N} \sum_{i=1}^N \left( y^{(i)} - w^T \hat{x}^{(i)} \right)^2$$

$$J(W) = \frac{1}{N} (Y - XW)^T (Y - XW)$$

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

$$J(W) = \frac{1}{N} (Y^T Y - 2Y^T XW + W^T X^T XW)$$

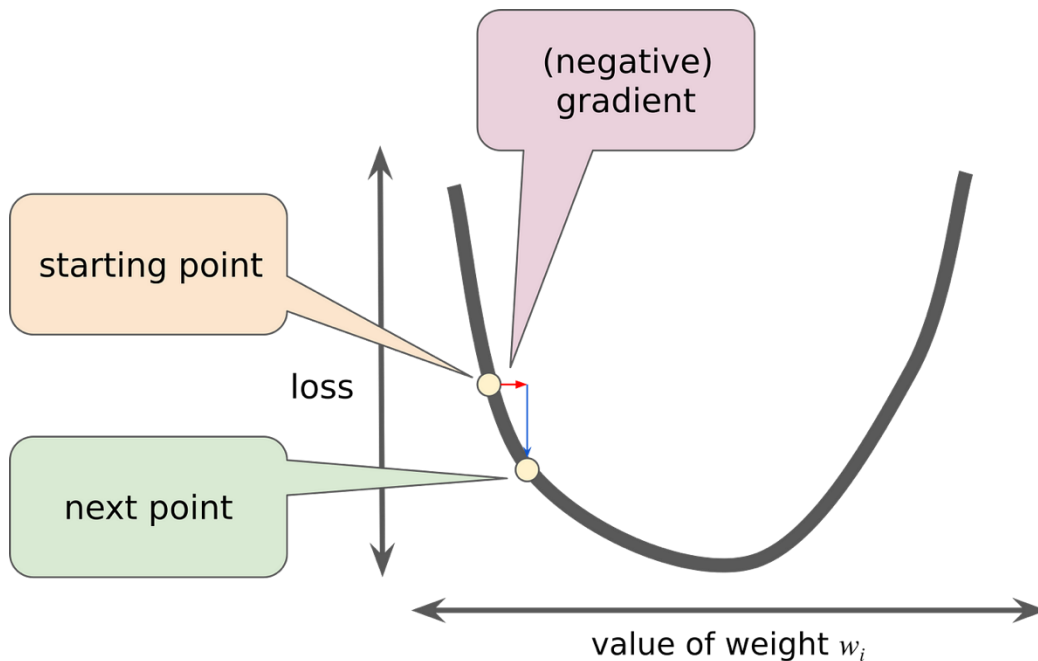
$$\nabla J(W) = \frac{\partial}{\partial W} \left[ \frac{1}{N} (Y^T Y - 2Y^T XW + W^T X^T XW) \right] \quad X^T XW = X^T Y$$

$$W^* = (X^T X)^{-1} X^T Y$$

$$\nabla J(W) = -\frac{2}{N} X^T Y + \frac{2}{N} X^T XW$$

# Linear Regression

$$J(w) = \frac{1}{N} \sum_{i=1}^N \left( y^{(i)} - w^T \hat{x}^{(i)} \right)^2$$



**If the gradient of a function is non-zero at a point, the direction of the gradient is the direction in which the function increases most quickly**

# Linear Regression

$$J(w) = \frac{1}{N} \sum_{i=1}^N \left( y^{(i)} - w^T \hat{x}^{(i)} \right)^2$$

$$\nabla J(w) = \frac{\partial J(w)}{\partial w}$$

$$\nabla J(w) = -\frac{2}{N} \sum_{i=1}^N (y^{(i)} - w^T \hat{x}^{(i)}) \hat{x}^{(i)}$$

$$J(w) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - w^T \hat{x}^{(i)})^2$$

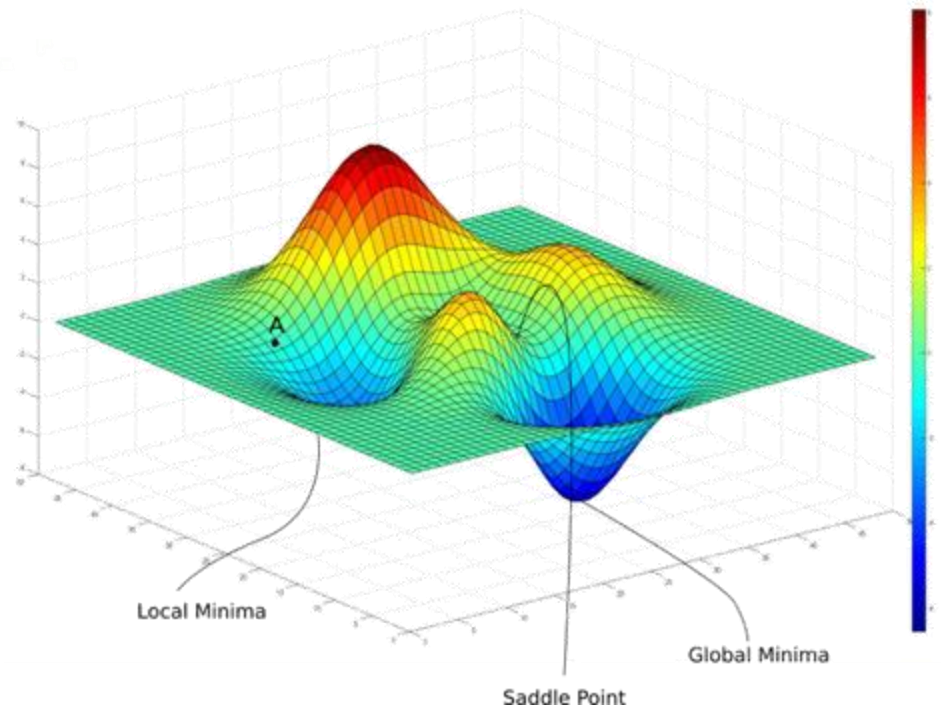
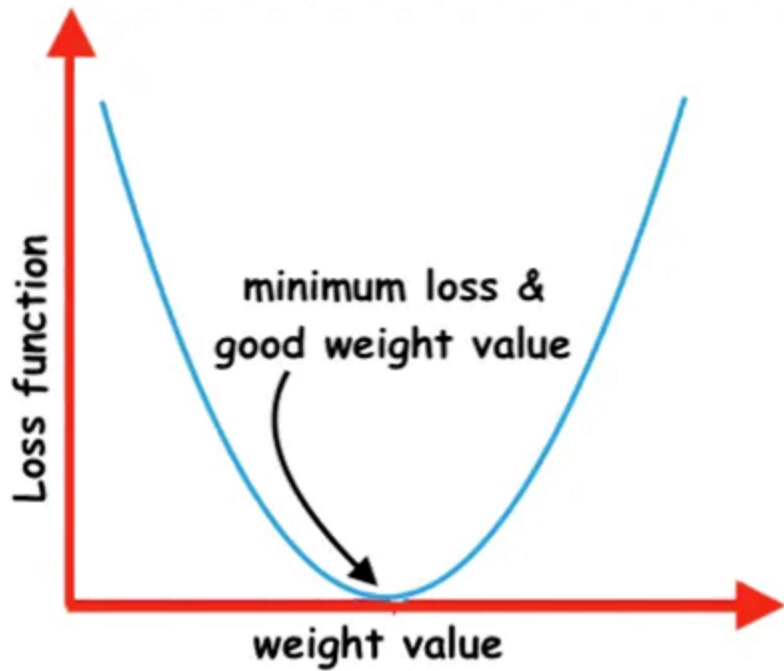
$$w_0 := w_0 + \frac{2\alpha}{N} \sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x_1^{(i)}))$$

$$\frac{\partial J(w)}{\partial w} = \frac{1}{N} \sum_{i=1}^N 2(y^{(i)} - w^T \hat{x}^{(i)}) (-\hat{x}^{(i)})$$

$$w_1 := w_1 + \frac{2\alpha}{N} \sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x_1^{(i)})) x_1^{(i)}$$

$$= -\frac{2}{N} \sum_{i=1}^N (y^{(i)} - w^T \hat{x}^{(i)}) \hat{x}^{(i)}$$

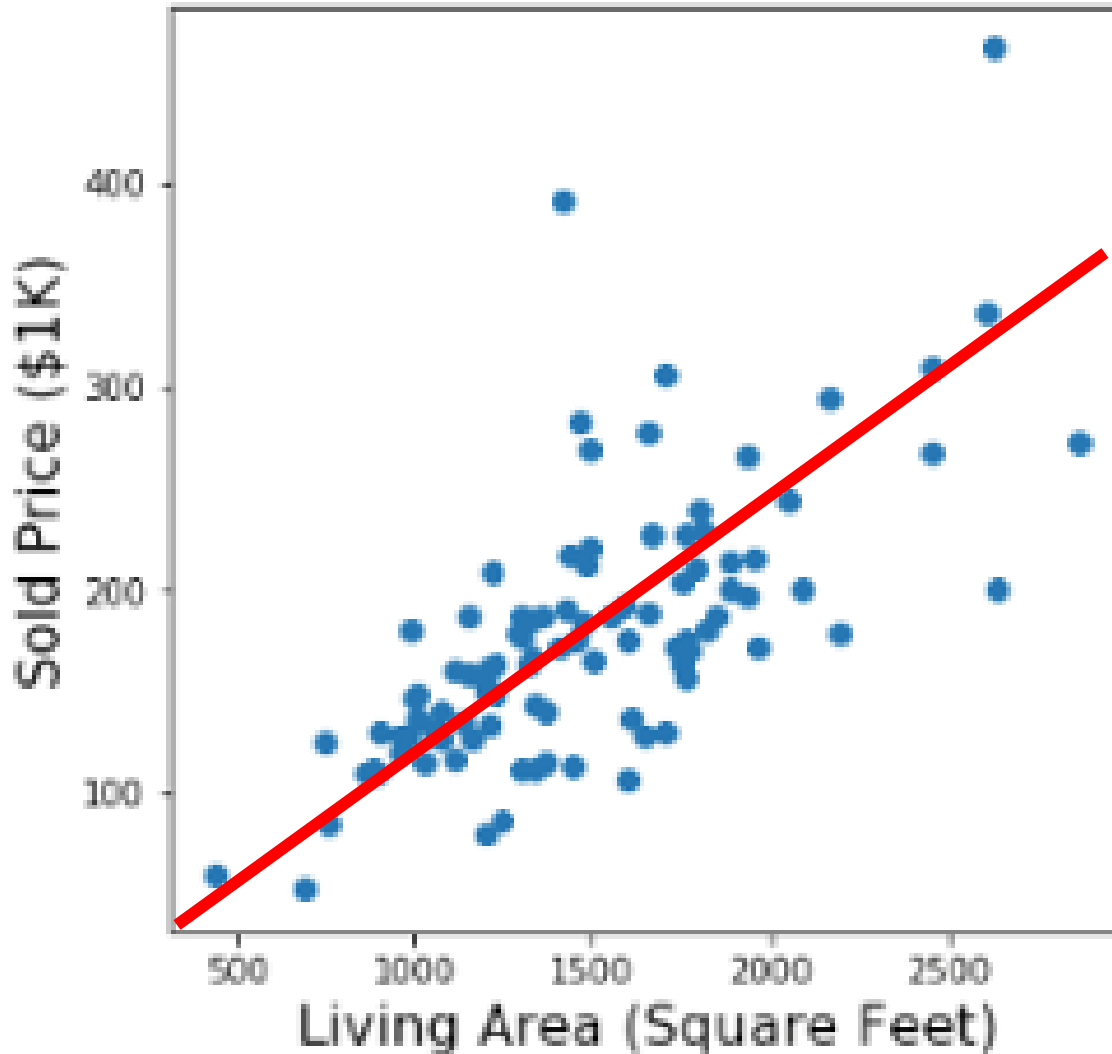
# Local Minimum vs Global Minimum



# Linear Regression

---

---

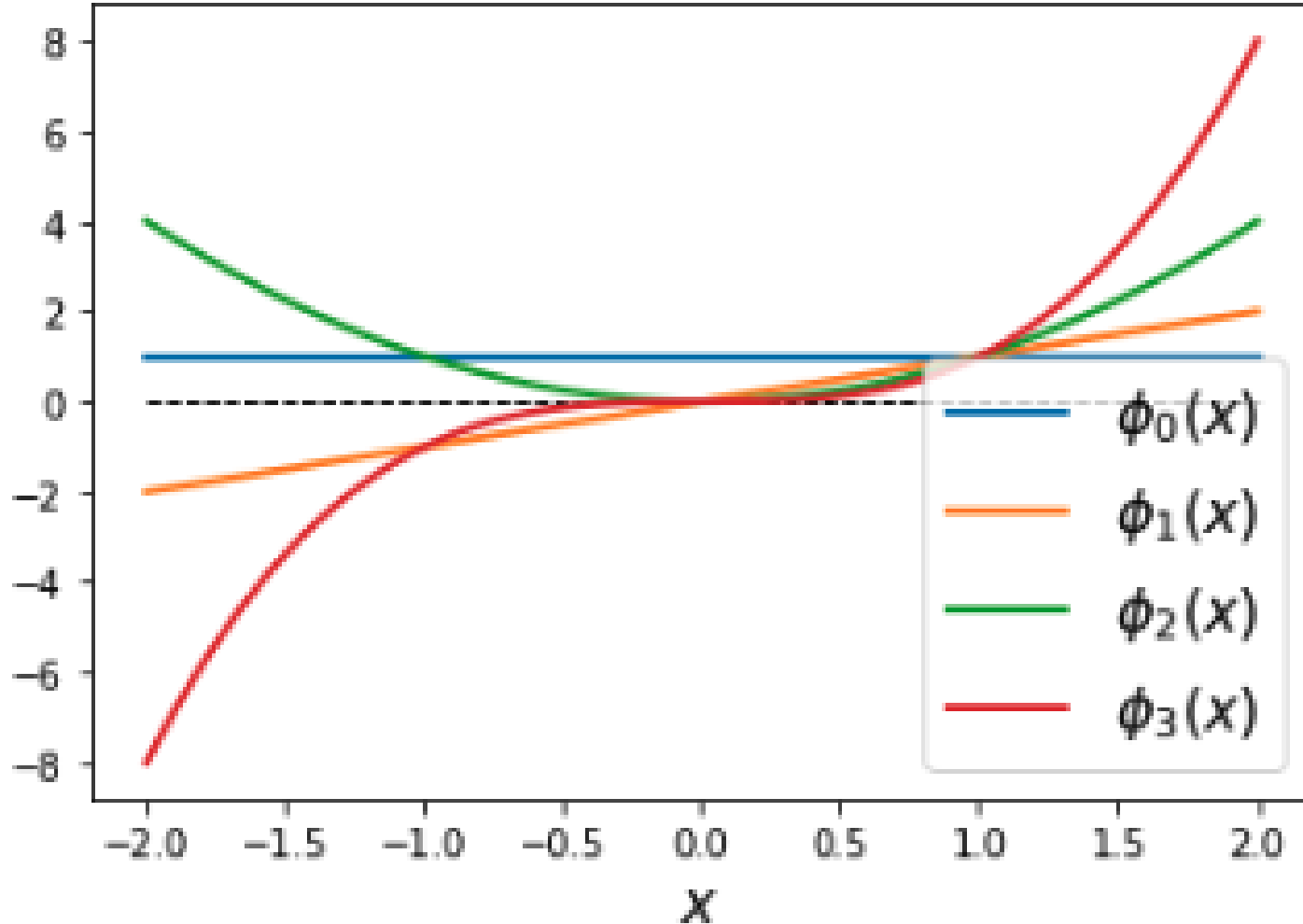


$$f(x) = w_0x + b$$

**Slope**

**Intercept**

# Problem of Linear Regression



# NonLinear Regression

---

- So far, we have been using a linear function for regression:

$$f(x) = w^T x + w_0 = \sum_{i=0}^d w_i x_i \quad (\text{Assuming } x_0 = 1)$$

- Lets generalize this model:

$$f(x) = \sum_{i=0}^M w_i \phi_i(x) = w^T \phi(x)$$

where  $\phi_i$  are fixed “basis” functions.

- For linear regression  $M = d$ ,  $\phi_i(x) = x_i$ .



# NonLinear Regression

$$f(x) = \sum_{i=0}^M w_i \phi_i(x) = w^T \phi(x)$$

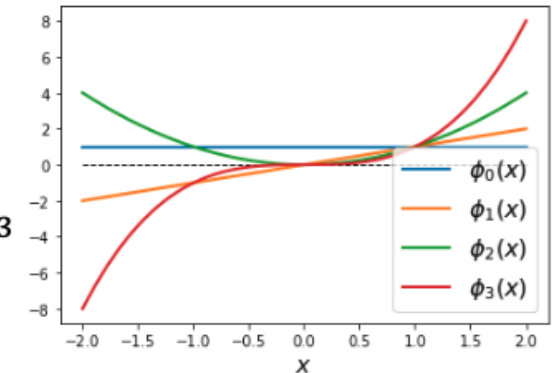
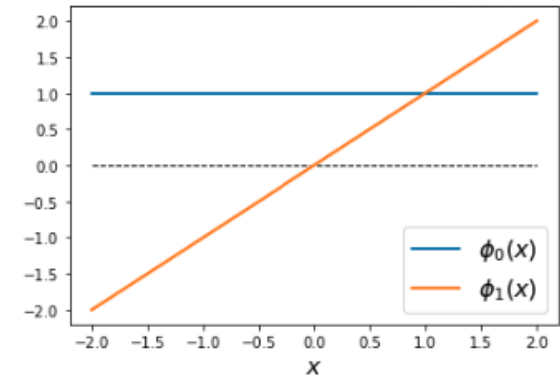
**E.g., Polynomial Regression:**

- 1D Polynomial Regression,  $\phi(x) = [1, x, x^2, x^3]$ :

$$\operatorname{argmin}_w \frac{1}{N} \sum_{n=1}^N (w^T \phi(x^{(i)}) - y^{(i)})^2$$

To avoid confusion, note that:  $\phi(x^{(i)}) = [1, x^{(i)}, (x^{(i)})^2, (x^{(i)})^3]$

$$f(x^{(i)}) = w_0 + w_1 x^{(i)} + w_2 (x^{(i)})^2 + w_3 (x^{(i)})^3$$



# NonLinear Regression

**Loss:** 
$$\operatorname{argmin}_w \frac{1}{N} \sum_{n=1}^N (w^T \phi(x^{(i)}) - y^{(i)})^2 = \operatorname{argmin}_w \|\Phi w - y\|^2$$

Where  $\Phi = [\phi(x^{(1)}), \dots, \phi(x^{(N)})]^T \in \mathbb{R}^{N \times M}$  and  $w \in \mathbb{R}^M$ .

## Optimization:

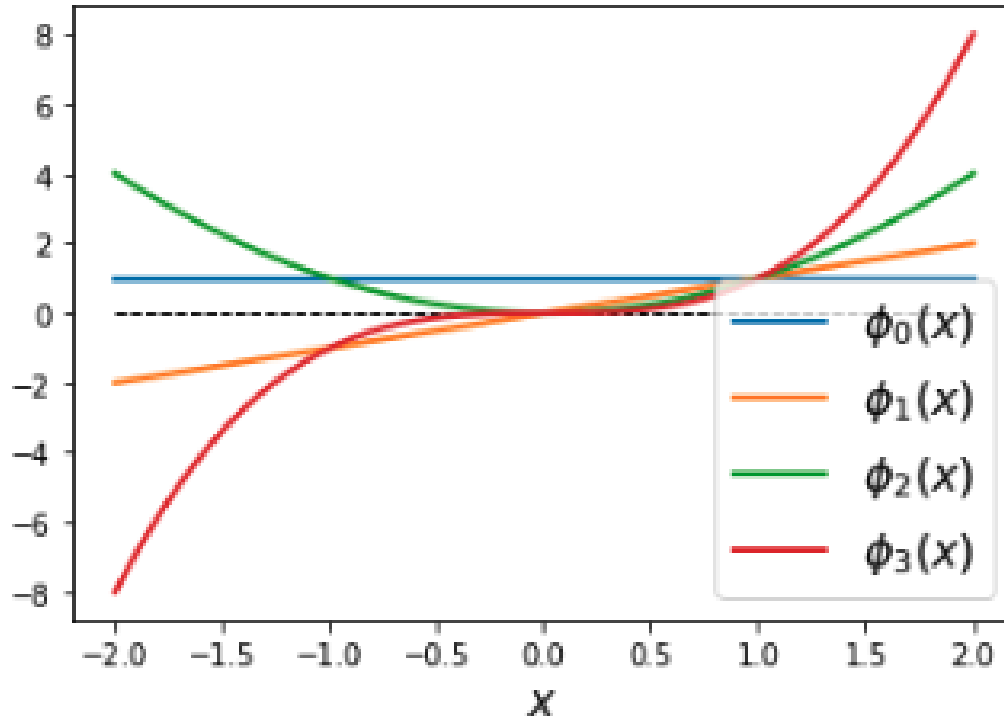
1. Closed form solution:  $w^* = (\Phi^T \Phi)^{-1} \Phi^T y$
2. Gradient descent:  $w^{(t)} = w^{(t-1)} - \epsilon \nabla_w \text{Loss}(w^{(t-1)})$

$$J(w) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - w^T \hat{x}^{(i)})^2$$

$$W^* = (X^T X)^{-1} X^T Y$$

**What is the problem of  
Nonlinear Regression?**

# NonLinear Regression



**The basis function is all fixed!**

**Can we learn the basis function?**

# NonLinear Regression

---

- Lets first look at what the learning problem might look like:

$$\operatorname{argmin}_w \sum_i \left( \left( \sum_j w_j \phi_j(x^{(i)}) \right) - y^{(i)} \right)^2$$

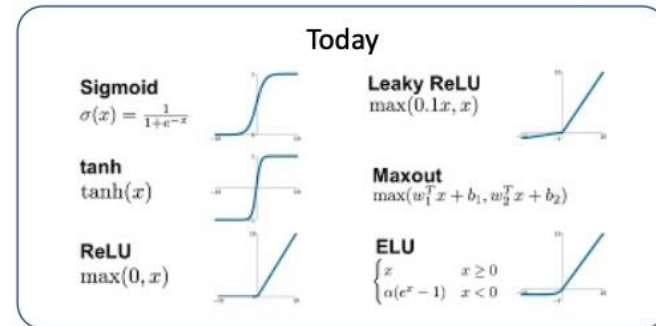
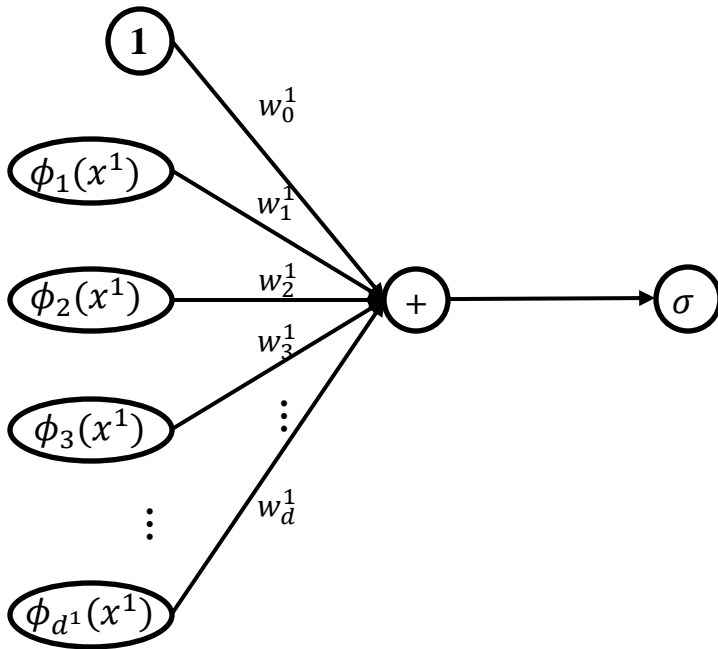
Neural Networks do this for us!

**What things are learned here?**

**What things are fixed here?**

# NonLinear Regression

## 1-layer Multi-layer Perceptron



$$\sum_{j=0}^{d^1} w_j^1 \phi_j(x^1)$$

# NonLinear Regression

---

- Lets first look at what the learning problem might look like:

$$\operatorname{argmin}_{w, \{\phi_j\}_{j=1}^M} \sum_i \left( \left( \sum_j w_j \phi_j(x^{(i)}) \right) - y^{(i)} \right)^2$$

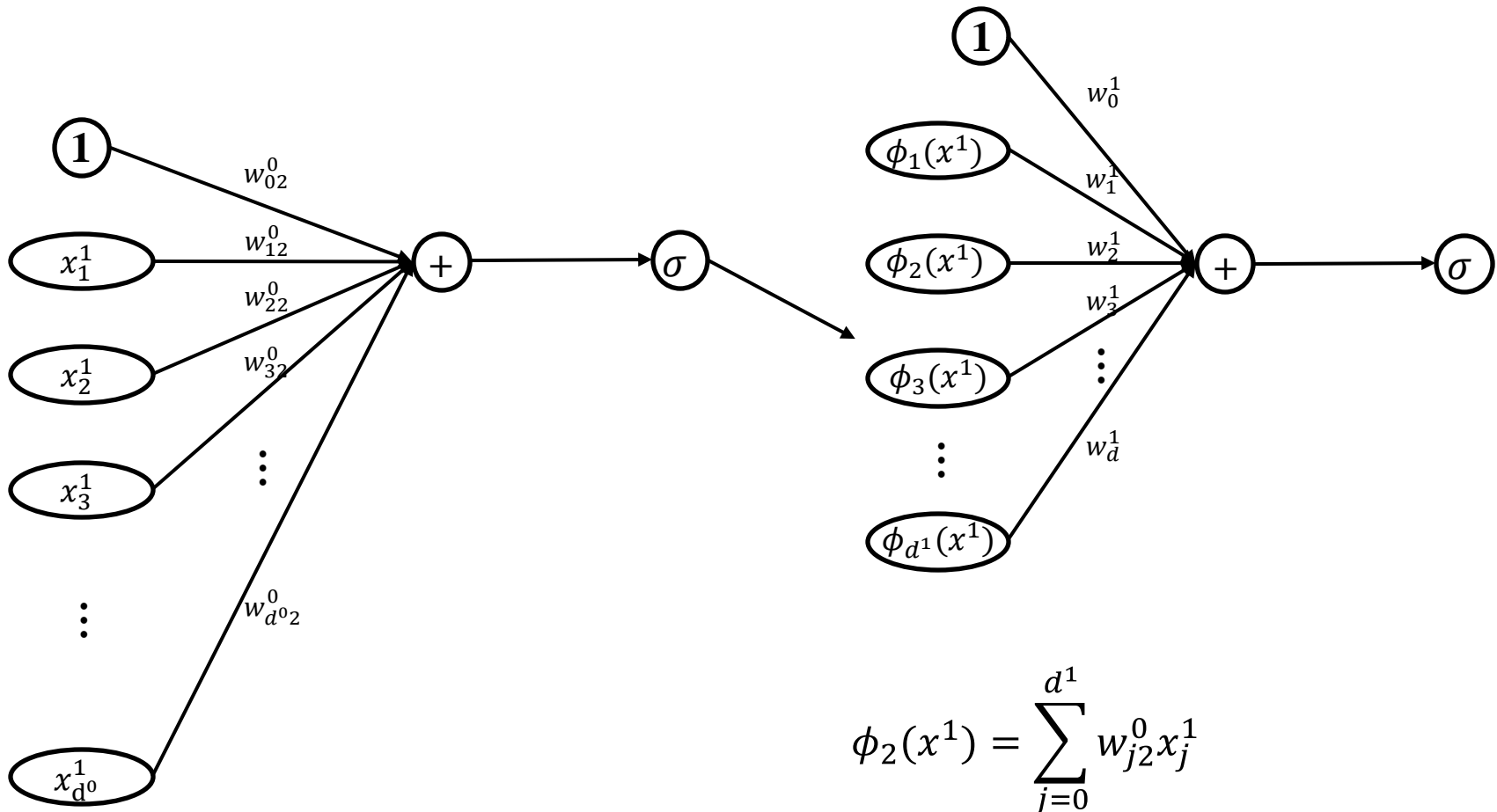
Neural Networks do this for us!

**What things are learned here?**

**What things are fixed here?**

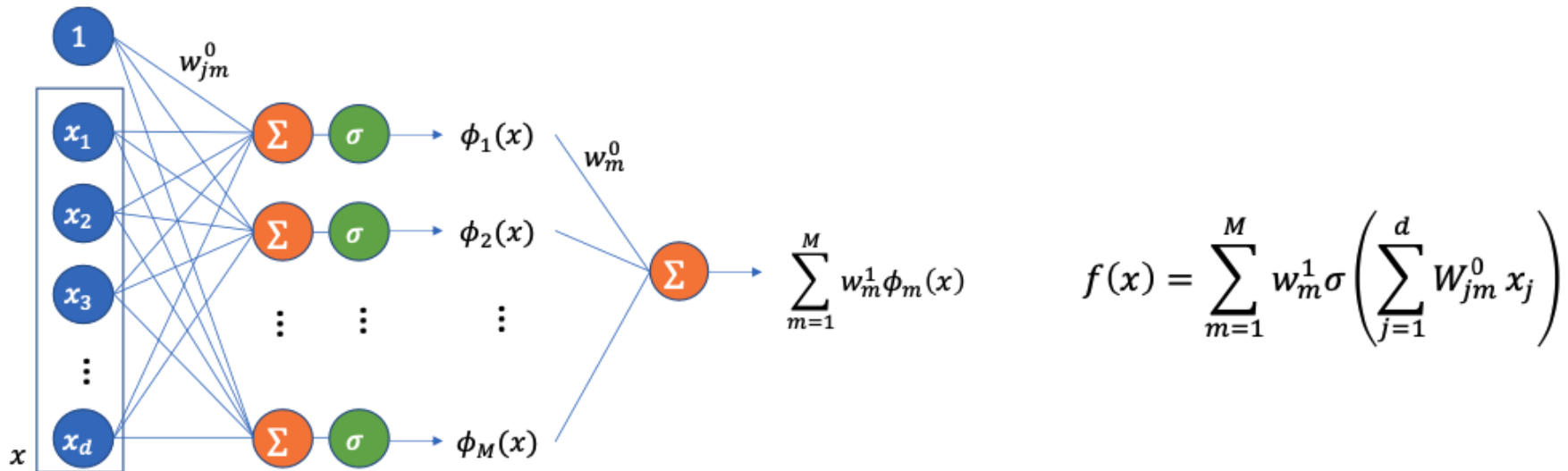
# NonLinear Regression

## 1-hidden layer Multi-layer Perceptron



# NonLinear Regression

## 1-hidden layer Multi-layer Perceptron

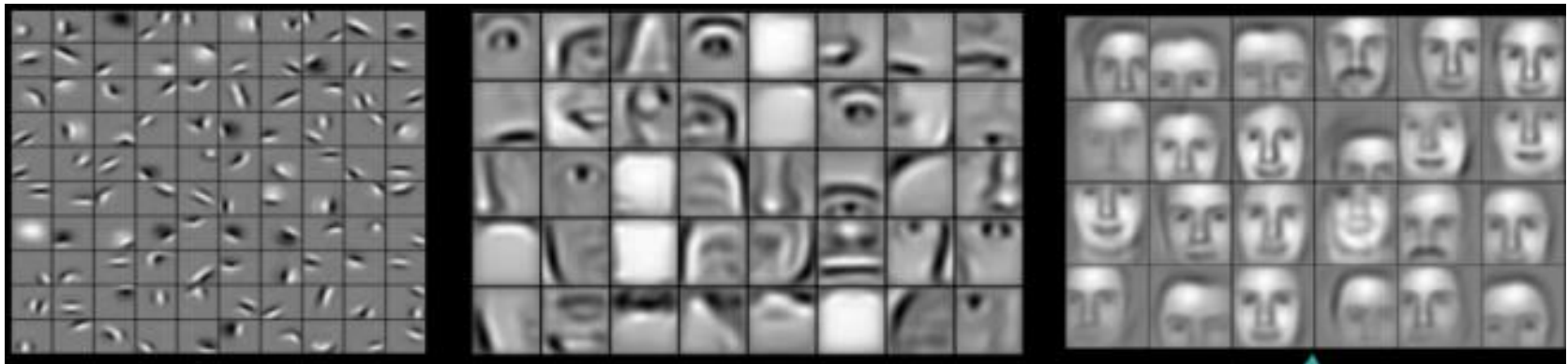




# Example

---

- Activations at hidden layers can be viewed as features extracted as functions of inputs
- Every hidden layer represents a level of abstraction
  - *Complex features are compositions of simpler features*



- Number of layers is known as depth of ANN
  - *Deeper networks express complex hierarchy of features*

# Question?

---



1. "Judge a man by his questions rather than by his answers."  
– Voltaire
2. "If I had an hour to solve a problem, I'd spend 55 minutes thinking about the problem and 5 minutes thinking about solutions."  
– Albert Einstein
3. "The art and science of asking questions is the source of all knowledge."  
– Thomas Berger
4. "Asking the right questions takes as much skill as giving the right answers."  
– Robert Half
5. "The wise man doesn't give the right answers, he poses the right questions."  
– Claude Lévi-Strauss
6. "Great questions make great companies."  
– Peter Drucker