# Intro to Value function approximation

Sunday, May 10, 2020    4:00 PM

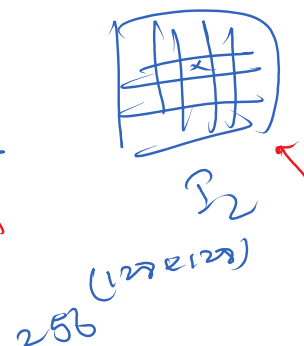<div style="text-align:center">

## Reinforcement Learning
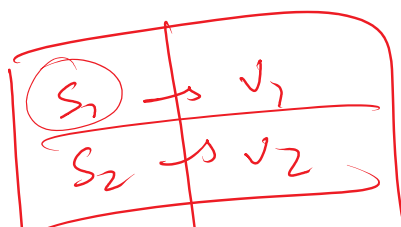### Value function approximation

</div>

**High-dimension state space:**



$10^{45}$

128

128

$S_1$

$0 - 255$

$S_2$

$256 (128 \times 128)$

*Function approximation*

$S_1 \to v_1$

$S_2 \to v_2$

$$S_2 \to v_2$$

**Supervised-learning perspective for approximation:**

*Monte − Carlo update equation:*
$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

*TD(0) update equation:*
$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

$$\text{New } \mathcal{EA} \leftarrow \text{Old } \mathcal{EA} + \alpha\left[\text{Target} - \text{Old } \mathcal{EA}\right]$$
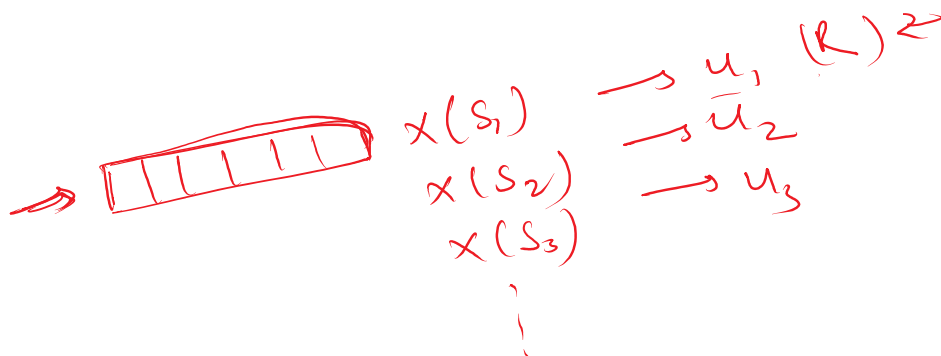
$$S_t \longrightarrow U_t$$

$$S_1 \longrightarrow U_1$$
$$S_2 \longrightarrow U_2$$
$$\vdots$$
$$S_T \longrightarrow U_T$$

$$X \longrightarrow Y$$

$$x \longrightarrow$$

$$\hat{f}(x) \longrightarrow \hat{y}$$

$$x(S_1) \longrightarrow U_1, (R)^2$$
$$x(S_2) \longrightarrow \bar{U}_2$$
$$x(S_3) \longrightarrow U_3$$
$$\vdots$$

$$x(S_t) \longrightarrow \boxed{\hat{v}(s, w)} \longrightarrow \hat{v}(S_t, w)$$

# Value Error

Reinforcement Learning
Value Error

**Value Error**

$$\mu(s)$$

$$\sum_s \mu(s) = 1$$

$$(\hat{y}_i - y_i)^2$$

$$\left[ \hat{v}(s,w) - v_\pi(s) \right]^2$$

*Mean squared Value Error is defined as:*

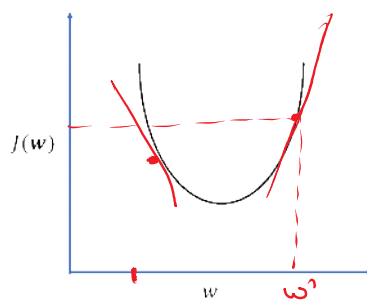$$\overline{VE}(w) = \sum_{s \in S} \mu(s) \left[ v_\pi(s) - \hat{v}(s,w) \right]^2$$

$$w^* \longrightarrow VE(w^*) \leq VE(w)$$

# Gradient based learning

## Reinforcement Learning
### Gradient Monte-Carlo & Semi-gradient TD

**Stochastic Gradient Descent**



$$J(w) = \sum [y_i - f(w)]^2$$

$$\nabla J(w) = 2 * \sum [y_i - f(w)] \nabla f(w)$$

$$w \leftarrow w - \alpha \nabla J(w)$$

$$w' \geq w \quad \text{if} \quad \nabla \geq 0$$

## Stochastic Gradient Descent on Value Error

$$\tilde{v}(s,w)$$

$$S_t \qquad u_t$$

$$\mu(s)$$

*SGD weight update for Value Error:*

$$w_{t+1} = w_t - \frac{1}{2}\alpha\nabla[v_\pi(S_t) - \hat{v}(S_t,w_t)]^2$$
$$= w_t + \alpha[v_\pi(S_t) - \hat{v}(S_t,w_t)]\nabla\hat{v}(S_t,w_t)$$

$$\nabla\hat{v}(S_t,w_t) = \left(\frac{\partial\nabla\hat{v}(S_t,w_t)}{\partial w_t^{(1)}}, \frac{\partial\nabla\hat{v}(S_t,w_t)}{\partial w_t^{(2)}}, \frac{\partial\nabla\hat{v}(S_t,w_t)}{\partial w_t^{(3)}}, \dots, \frac{\partial\nabla\hat{v}(S_t,w_t)}{\partial w_t^{(d)}}\right)$$

## Gradient Monte-Carlo

*Gradient Monte − carlo weight update rule:*

$$w_{t+1} = w_t + \alpha[G_t - \hat{v}(S_t,w_t)]\nabla\hat{v}(S_t,w_t)$$

---

**Gradient Monte Carlo Algorithm for Estimating $\hat{v} \approx v_\pi$**

Input: the policy $\pi$ to be evaluated
Input: a differentiable function $\hat{v} : S \times \mathbb{R}^d \to \mathbb{R}$
Algorithm parameter: step size $\alpha > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop forever (for each episode):
  Generate an episode $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$ using $\pi$
  Loop for each step of episode, $t = 0, 1, \dots, T-1$:
    $\mathbf{w} \leftarrow \mathbf{w} + \alpha\left[G_t - \hat{v}(S_t,\mathbf{w})\right]\nabla\hat{v}(S_t,\mathbf{w})$

## Semi-Gradient TD(0)

*Semi Gradient TD(0) weight update rule:*

$$w_{t+1} = w_t + \alpha[R_{t+1} + \gamma\hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t)]\nabla\hat{v}(S_t, w_t)$$

**Semi-gradient TD(0) for estimating** $\hat{v} \approx v_\pi$

Input: the policy $\pi$ to be evaluated
Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \to \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$
Algorithm parameter: step size $\alpha > 0$
Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        Choose $A \sim \pi(\cdot|S)$
        Take action $A$, observe $R, S'$
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R + \gamma\hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})]\nabla\hat{v}(S, \mathbf{w})$
        $S \leftarrow S'$
    until $S$ is terminal

# Linear Function Approximation

Reinforcement Learning
Linear Function Approximation

**State value function using linear function approximator**

$$S \longrightarrow x(S) \qquad |x(s)| = |w|$$

$$x(s) = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}$$
$$w = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix}^T$$

$$\hat{v}(s,w) = w^T x(s) = \sum_{i=1}^{d} w_i x_i(s)$$

$$\hat{v}(s,w) = \left( x_1 w_1 + x_2 w_2 + x_3 w_3 \right)$$

$$\nabla \hat{v}(s,w) = \left[ \frac{\partial \nabla \hat{v}(s,w)}{\partial w_1}, \frac{\partial \nabla \hat{v}(s,w)}{\partial w_2}, \frac{\partial \nabla \hat{v}(s,w)}{\partial w_3} \right]$$

$$\nabla \hat{v}(s,w) = x(s)$$
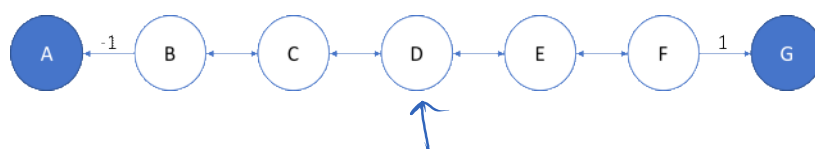
$$\nabla \hat{v}(s, w) = \boldsymbol{x}(s)$$

*SGD update for linear function approximator*:
$$w_{t+1} = w_t + \alpha[U_t - \hat{v}(S_t, w_t)]\boldsymbol{x}(s)$$

$$\frac{\partial v}{\partial w_3}$$

$$= [x_1 \quad x_2 \quad x_3]$$

$$= x(s)$$

$$\overline{VE}(w_{TD}) \leq \frac{1}{1-\gamma}\left(\overline{VE}(w_*)\right)$$

# Reinforcement Learning
## State Aggregation



$D: \quad x(D) = [ 0 \quad 0 \quad 0 \quad 1 \quad 0 \quad 0 \quad 0 ]$

$A: \quad x(A) = [ 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 ]$

$B: \quad x(B) = [ 0 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 ]$

$w = [ w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6 \quad w_7 ]$

$D: \quad x(D) \longrightarrow u_t$

*SGD update for linear function approximator:*
$$w_{t+1} = w_t + \alpha[U_t - \hat{v}(S_t, w_t)]x(s)$$

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \\ w_7 \end{bmatrix} \quad 2 \quad \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \\ w_5 \\ w_6 \\ w_7 \end{bmatrix} \quad f \quad C \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$10^{n8}$   States   $10^6$

$(1 - 100) : 1$
$[01 - 200 : 2]$   $10^6$   $[1 \ 0 \ 0 \ 0 \ 0 \ \cdots \ 0]$
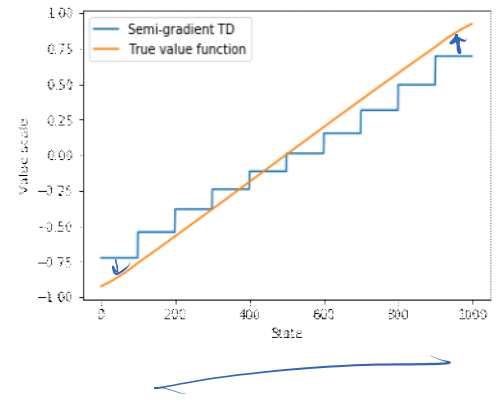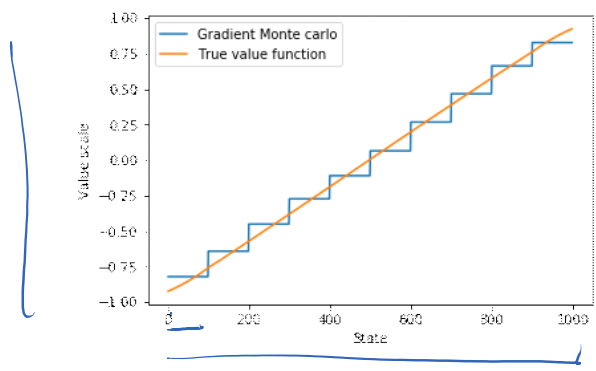$[0 \ 1 \ 0 \ 0 \ \cdots ]$
$\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

## 1000 state random-walk MDP with state aggregation:

- States are numbered form 1 to 1000, from left to right.
- All episode begin near the center on state 500.
- States 1 and 1000 are terminal states, and it fetches a reward of -1 and 1 on transitioning to them, respectively.
- Policy:
  - Agent can transition from current state to one of 100 states in left or rights, all with equal probability. So there are 200 possible actions. At edges, if an action goes beyond 1 or 1000, then the agent instead transition to the terminal states.
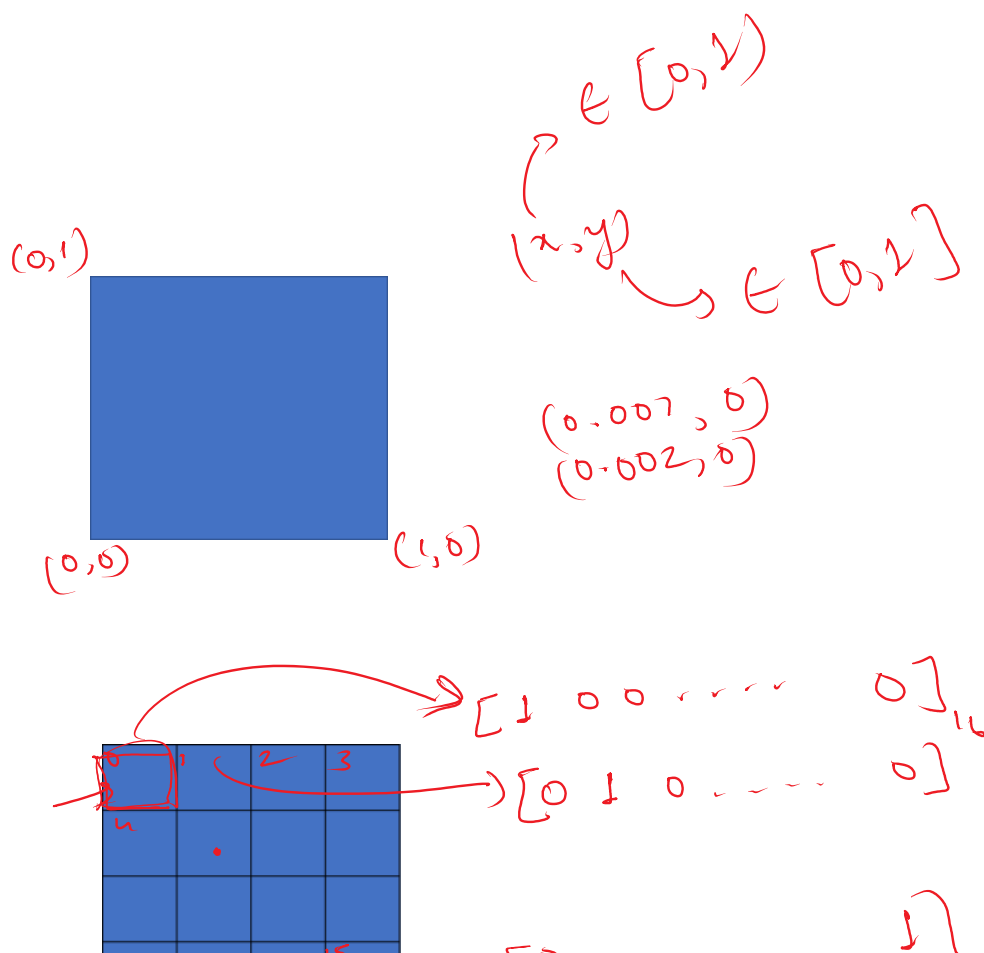


400   499  500   501   600

$[1 \ 00 \ 00000000 \ 0]$

$\underline{1 - 100} \rightarrow \underline{1} \quad ; \quad [1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$

$101 \rightarrow 200 \rightarrow 2 \quad [0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$

$201 \rightarrow 300 \rightarrow 3$

$901 \longrightarrow 1000 \rightarrow 10 \quad [0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1]$

# Reinforcement Learning
## Coarse coding and tile coding

**Coarse coding:**

$$\begin{pmatrix} x, y \end{pmatrix} \begin{cases} \in [0,1) \\ \in [0,1] \end{cases}$$

(0,1)

(0,0)          (1,0)

$(0.007, 0)$
$(0.002, 0)$

$[1\ 0\ 0\ \cdots\cdots\ 0]_{16}$
$[0\ 1\ 0\ \cdots\cdots\ 0]$
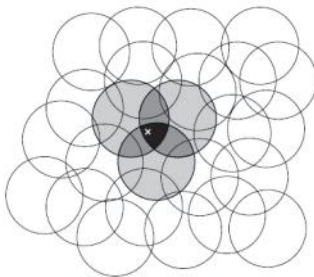
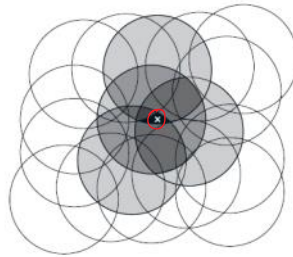$\rightarrow [0 \cdots$

n circles

$[0 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1]_n$
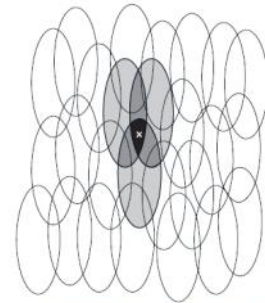
$\downarrow$

$[w_1 \ w_2 \cdots \qquad w_n]$



Narrow generalization    Broad generalization    Asymmetric generalization

**Tile coding:**



Tiling 1

Tiling 2

Tiling 1
Tiling 2
Tiling 3
Tiling 4

Continuous 2D state space

Point in state space to be represented

Four active tiles/features overlap the point and are used to represent it
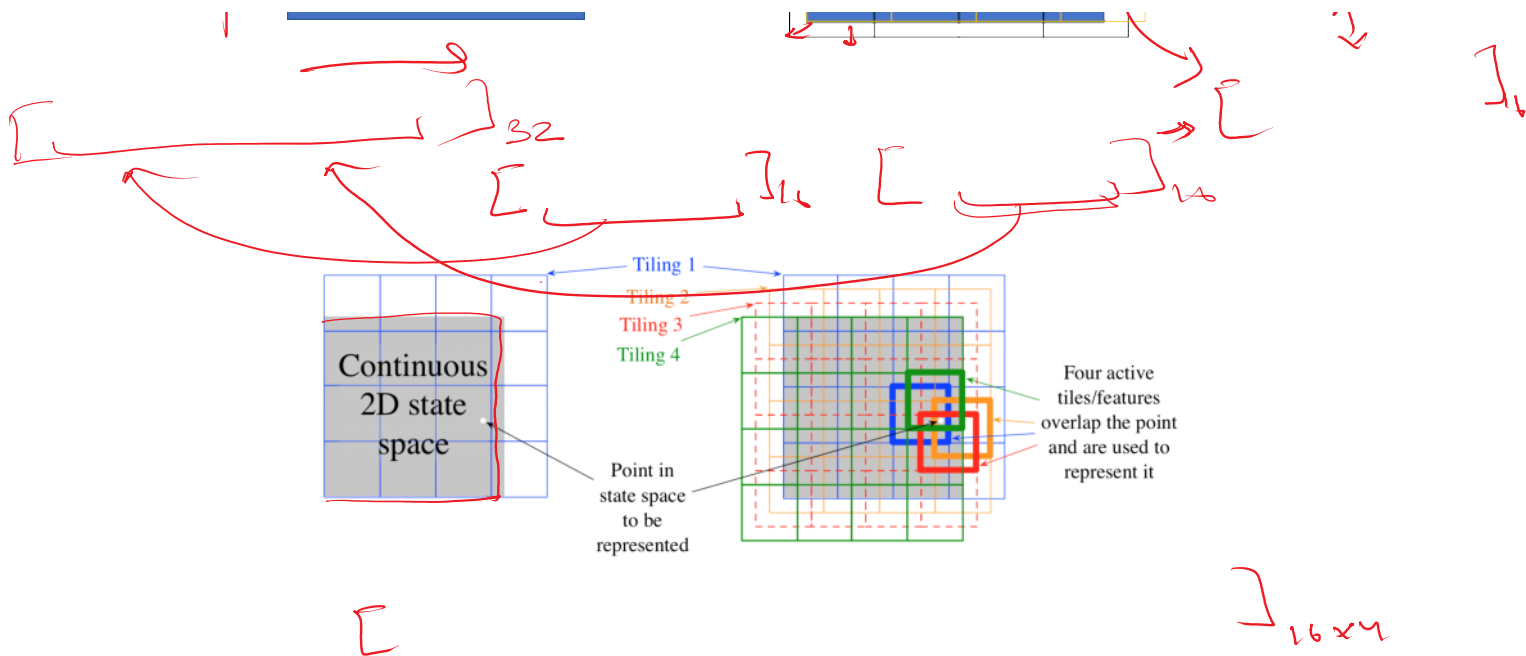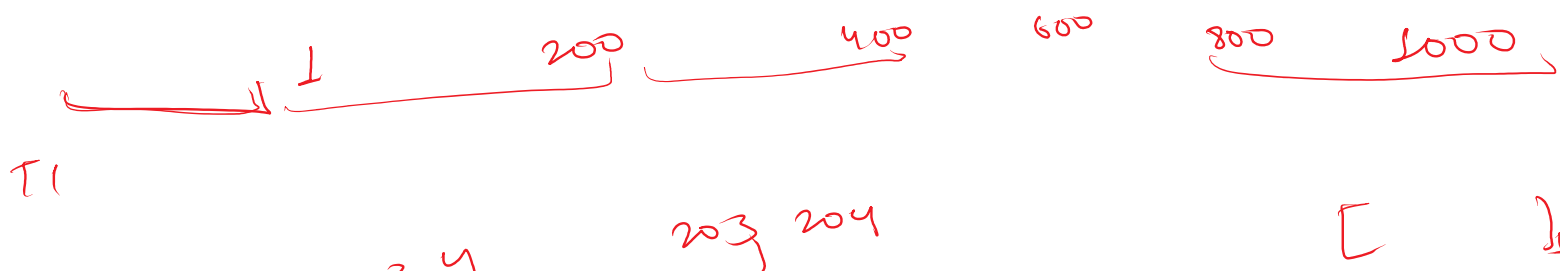
## 1000 state random-walk MDP:

- States are numbered form 1 to 1000, from left to right.
- All episode begin near the center on state 500.
- States 1 and 1000 are terminal states, and it fetches a reward of -1 and 1 on transitioning to them, respectively.
- Policy:
  - Agent can transition from current state to one of 100 states in left or rights, all with equal probability. So there are 200 possible actions. At edges, if an action goes beyond 1 or 1000, then the agent instead transition to the terminal states.
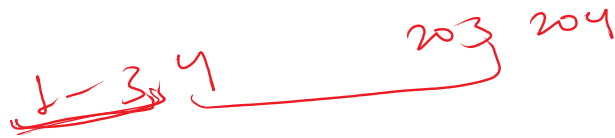
-1

1

| 1 | 2 | 3 | ........ | 998 | 999 | 1000 |

⌐⌐

T2    ↓−3↓4 _____ 203 204

[     ]

[     ]

T≤50

[          ]₆ ×50
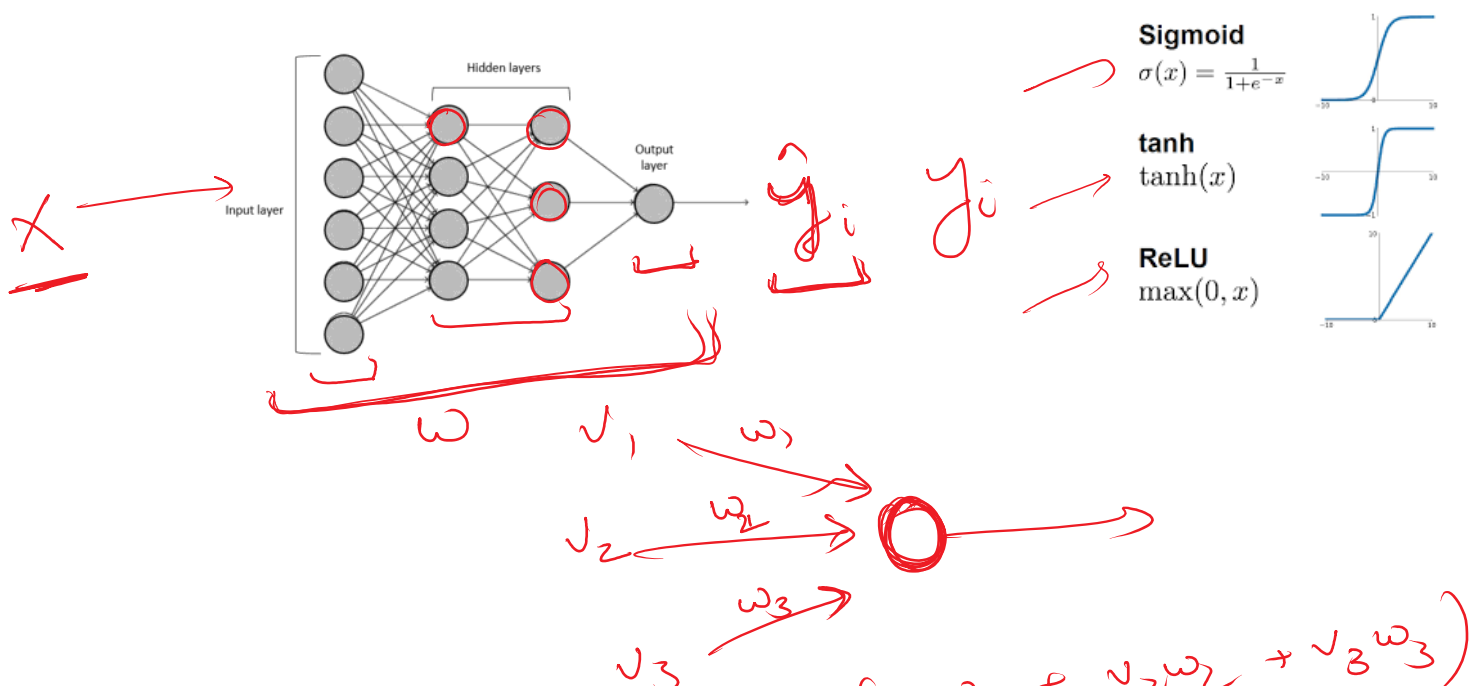
⌐→ 50 values = 1

[          ]₆ₓ₅₀

## Reinforcement Learning
### Non-linear function approximator - Artificial Neural Networks

**Artificial Neural Networks:**



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

**tanh**
$\tanh(x)$

**ReLU**
$\max(0, x)$

$$v_3 \xrightarrow{\omega_3} \quad \not{E}\left( v_1 \omega_1 + v_2 \omega_2 + v_3 \omega_3 \right)$$

$$x(s_t) \longrightarrow \boxed{NN(\omega)} \longrightarrow \hat{v}(s_k, \omega)$$