# ML week 11 - Reinforcement Learning

RF learning is loosely based on how humans learning

RL Agent learns by interacting w/ environment
↳ to achieve given goal(s)

## RL environment

State ➡ Agent observe the env state/ conditions
$s_t \in S$ @ time stamp $t$

State has a markov property. dependant on state @ tag of $t$'s. e.g $(t)$, $(t, t-1)$ $(t, t-1, t-2)$

Action ➡ Agent can select an action @ time $t$ ➡ $a_t$
➡ availble given the state
➡ $A(s_t)$

$$a_t \in A(s_t)$$

reward ➡ the reward that the env gives at $t+1$ dependant on the action taken @ $t$

Arrave @ new state $s_{t+1}$

▶ <u>RL Agent elements</u>

Agent ➤ · continually interacts w/ timesteps
· obj is to max reward

Policy ➤ Function that specifis the
Behaviours available to an agent
given the env

     State to a chon

     function may be either deterministic
or stochastic

value ➤ agents estimation of future reward


▶ <u>reward vs value</u>

reward is short-term gain from next ~~reward~~

value long-term total gain
    ↳ more important to learning

▶ <u>Supervised vs Reinforcement</u>

| SL | RL |
|---|---|
| • Model (weights) | • Agent (Policy) |
| • Input data | • Env state |
| • Prediction | • Agent Action, next state |
| • loss func | • Reward |
| Model is train on existing ground truth (input data) | Agent & Env interact to learn |
| min loss on train data, generalize to unseen | Agent learns from own experience (no exist truth) |
| | Max reward via Exploration of action |

▶ <u>Exploit vs Exploration</u>

Exploitation ⟶ Select largest reward for next
(greedy)      action, requires knowledge of
              action & rewards

Exploration ⟶ Select action w/ uncertain
              value to obtain better value
              estimate

              test & learn for future