



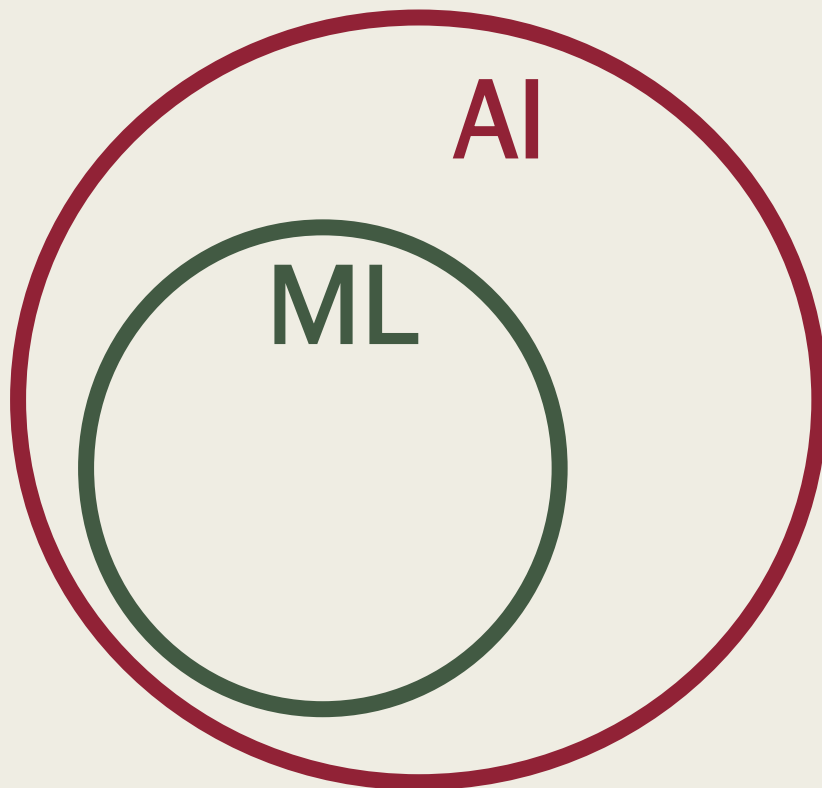
REINFORCEMENT LEARNING

torresjm – jan 2022





AI



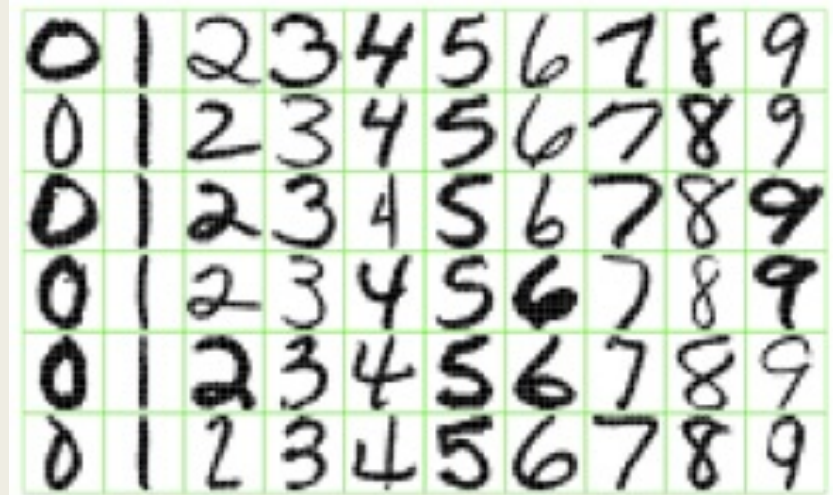
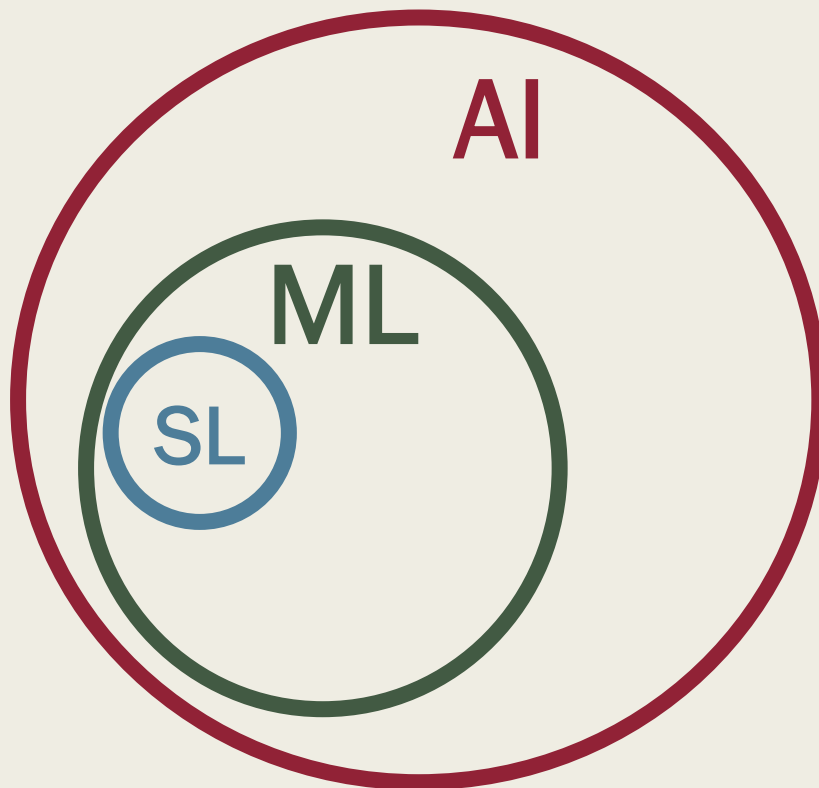
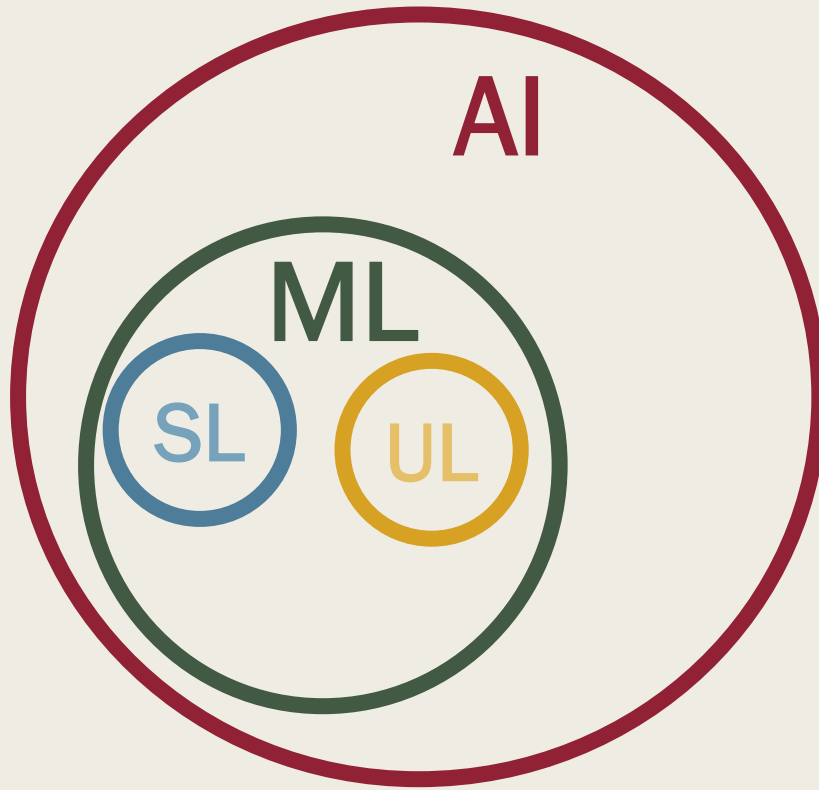


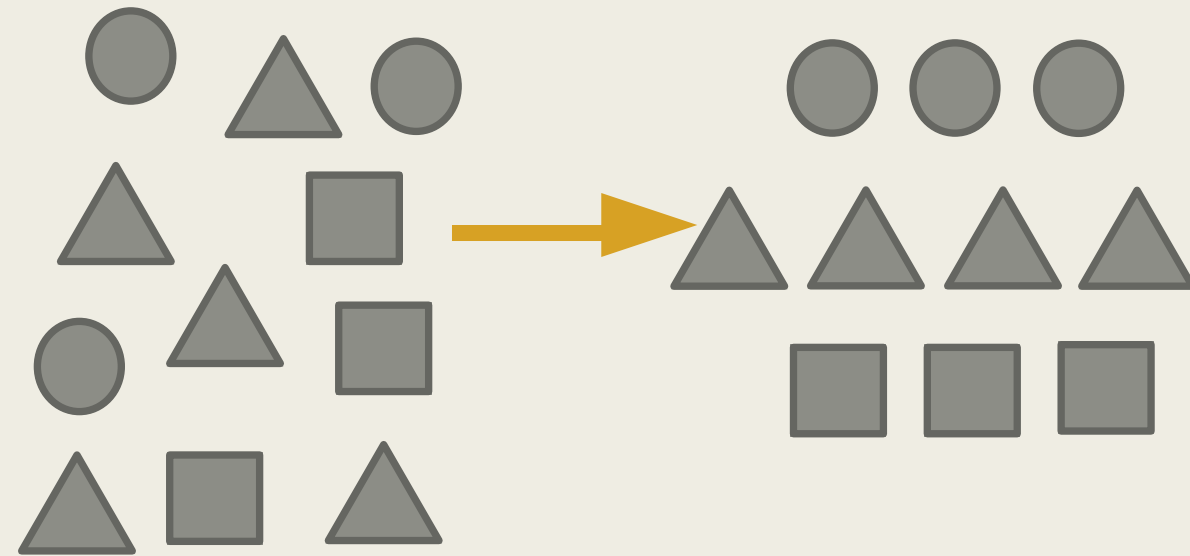
Figure 1.2: *Examples of handwritten digits from U.S. postal envelopes.*

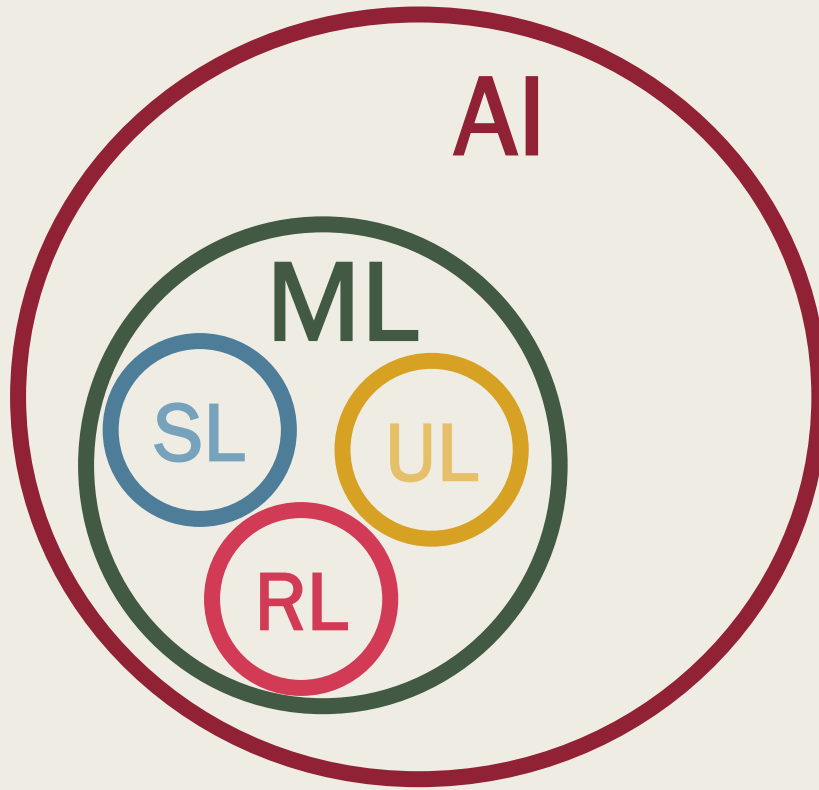
- Data, labels
- Train
- Test
- Predict



Data (and no label)

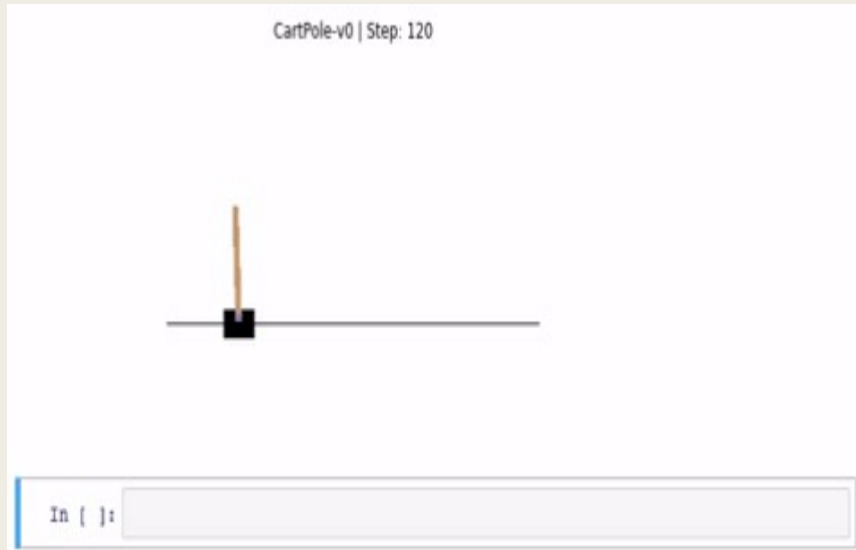
finding some structure in the dataset, grouping data.





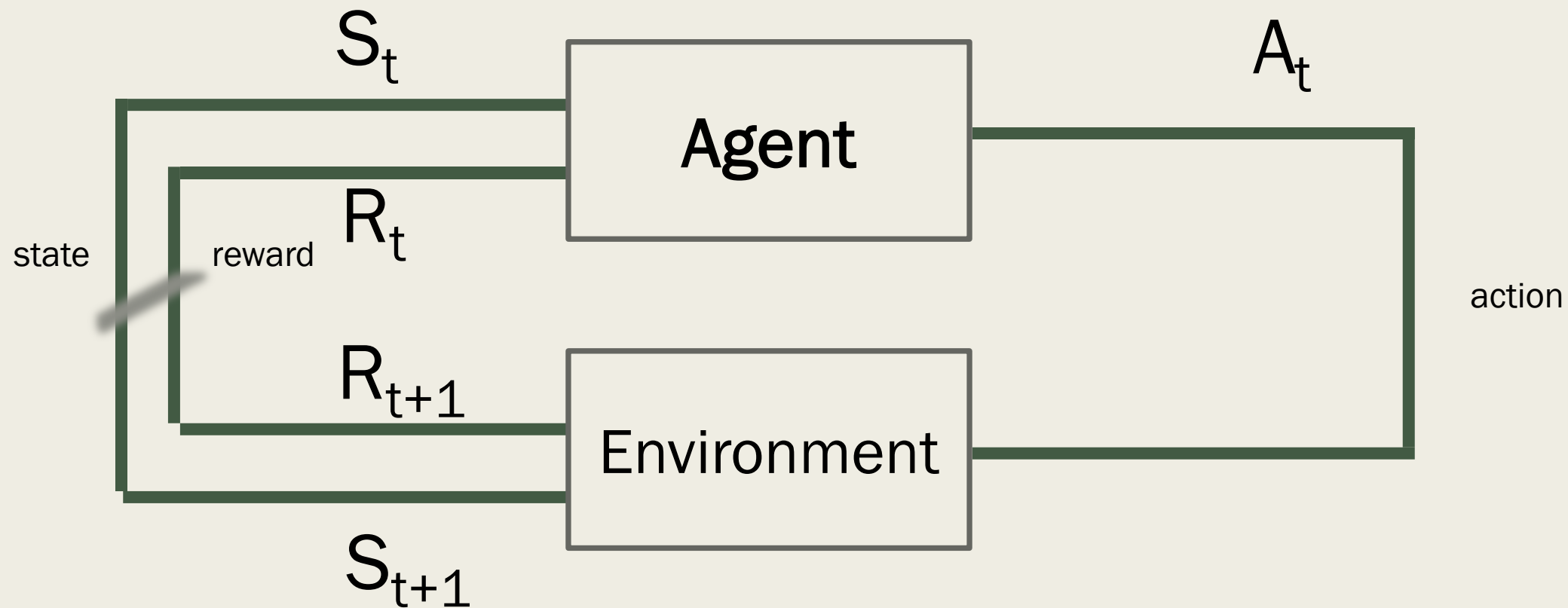
an « agent » interacts with an « environment » to learn what sequences of actions will maximize the rewards it will receive.

RL - example



Cartpole :

- the agent controls the cart, trying to keep the pole up. interaction ends when the pole falls down
- 2 actions :
 - left,
 - right
- State :
 - position,
 - cart velocity,
 - cart / pole angle,
 - top of pole velocity



trajectory : $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots R_n, S_n, A_n$

Decision process

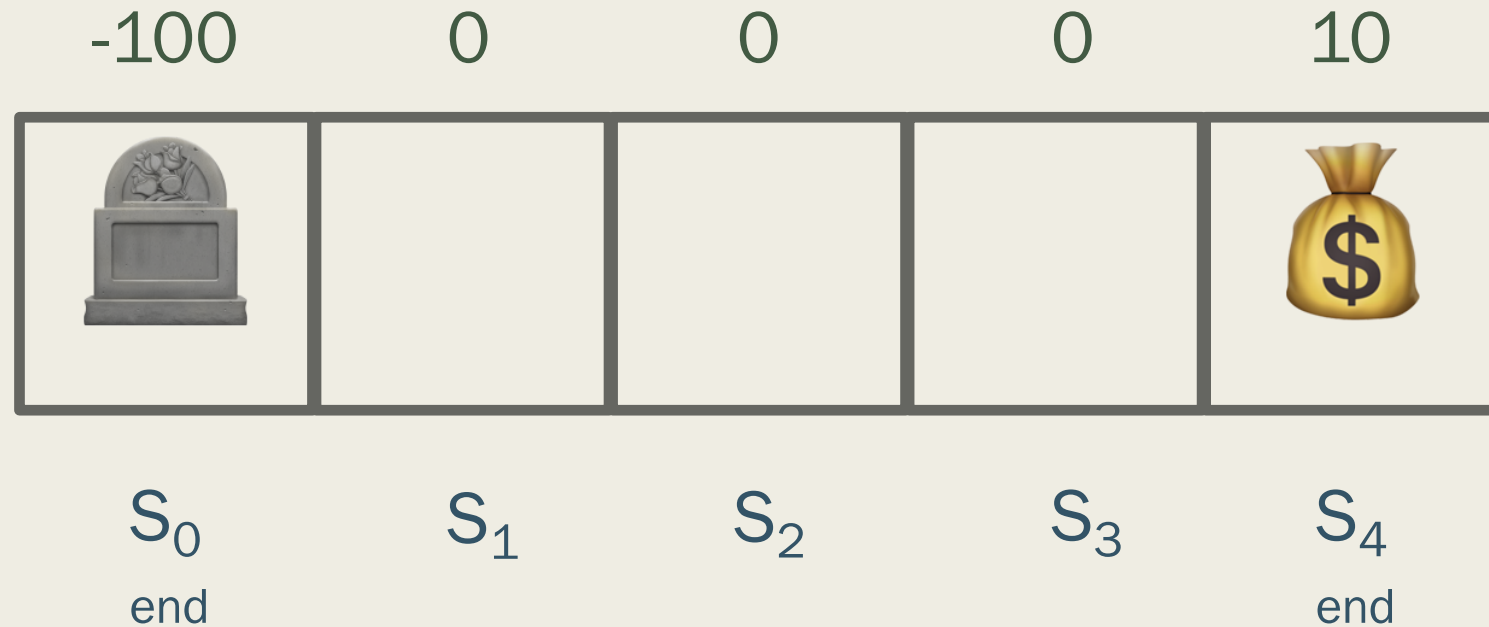
s, a, p, r

$$P(s' | s, a) = P(S_{t+1} | S_t = s, A_t = a)$$

Notes:

- markovian process (future depends on present, not past)
- p, r : are given (model based) or not (model free)

Example : Row GridWorld



Reward function : $r(s, a, s') = E[R_{t+1} \mid S_t = s, A_t = a, S_{t+1} = s']$

ex :

- $r(S_1, \leftarrow, S_0) = -100$
- $r(S_4, \leftarrow, S_3) = 0$
- $r(S_4, \rightarrow, S_5) = 10$

policy : π

- a policy π describes the behaviour of the agent (for each state in the environment)
- Resolving the environment mean finding the policy that maximizes the future return
- this will be the optimal policy π_*

Return is defined by : $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$

(more often : $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$)

note : $G_t = R_{t+1} + \gamma G_{t+1}$

For a given policy a state s defines an action a : $\pi(s) = a$

The state value function measures the return of a policy at each state :

$$v_{\pi}(s) = E_{\pi}[G_t \mid S_t = s]$$

policy : π

- Optimal policy π_* is associated to $v_*(s) = \max_{\pi} (v_{\pi}(s))$ for any s
- Now resolving the environment if finding v_*
- Resolving the environment is finding the policy that maximizes the future return

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma v_{\pi}(s')]$$

Which lead to the Bellman Equation (that give an algo for solving :

$$v_*(s) = \max_a [\sum_a \pi(a|s) \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma v_{\pi}(s')]]$$

try all policies,
finding v max
will give the
optimal policy

```
input  $\pi$ 
init  $v(s) = 0$  for all  $s$ 
repeat :
  for each  $s$  do :
     $v(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma v_{\pi}(s')]$ 
  end
until change of  $v$  is lower than a threshold
return  $v$  (=approx  $v_{\pi}$ )
```