# CS 3630 Reinforcement Learning: I

Markov Decision Processes

# Reinforcement Learning

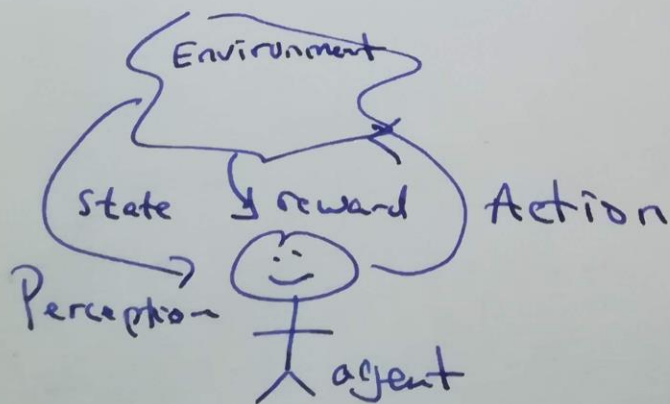## "Typical" Machine Learning

- Lots of Data
- passive Learner
- Learning = pattern analysis
  function approx.

### Examples

* Find photos of cats.
* product recommendation
* Text Translation

## Robots are Not Like This!



Environment

State ↓ reward   Action

Perception

agent

## RL:

RL is a process of modifying behavior by rewarding desired outcomes.

- Dogs — Treats, learn tricks
         — punishment
         → negative reward

- Kids — 1↯/A on report

- Students: grades
            joy, satisfaction
            starting Salary

1. Robot Senses world
2. Robot decides & executes action
3. World state changes
4. Robot Receives a reward. }
   Robot Senses world          }
5. Robot update its "strategy"
6. Go to 2

## Questions:

- Mathematical Model
  - states
  - actions
  - uncertainty (!)
- Mathematical formalism for reward
- Given the above How to compute an Optimal Strategy/policy
- Now... suppose we don't know the models for world, actions... How can we learn them?

---

- Markov Process
- Markov Decision Process
- Value function & how to compute it
- R.L.

## Markov Processes

Simplest case:

- discrete time, $k = 0, 1, 2, \ldots$
- discrete states, $S$ set of states.
- A set of transition probabilities

$$T : S \times S \longrightarrow \underbrace{[0,1]}_{\text{probability}}$$

$$P\{S_{k+1} \mid S_k\} = \underline{T}(s, s')$$
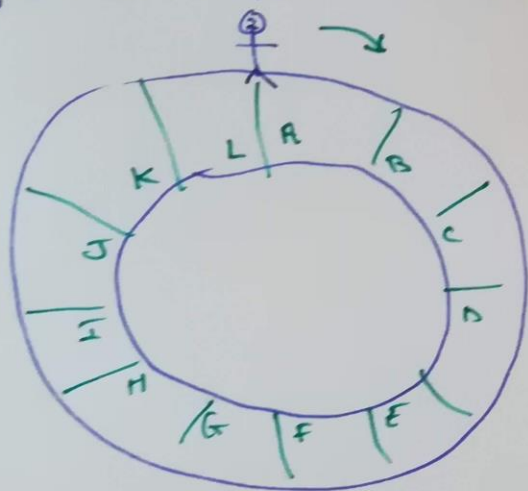
$$= T(S_k, S_{k+1})$$

$$T(s, s') = \text{prob }\{\text{arriving to state } s', \text{ given we are now in state } s\}$$

---

$\longrightarrow$ M.P. evolves "autonomously"

# Example:

Sisyphus has the job of, every day, throwing a large rock, on a circular track.



$L_k$ = distance of throw on $k^{th}$ day

$$P\{L_k = 1\} = 0.25$$

$$P\{L_k = 2\} = 0.5$$

$$P\{L_k = 3\} = 0.25$$

---

$$S = \{A, B, \ldots L\}$$

$$T(A,B) = 0.25 \quad T(A,C) = 0.5 \quad T(A,D) = 0.25$$

$$\vdots$$

$$T(J,K) = 0.25 \quad T(J,L) = 0.5 \quad T(J,A) = 0.25$$

$$T(K,L) = 0.25 \quad T(K,A) = 0.5 \quad T(K,B) = 0.25$$

$$T(s, s') = 0 \text{ otherwise}$$

---

## Properties:

- Stationary: T does not change over time

$$\boxed{P\{S_{k+1} \mid S_0, S_1, \ldots S_k\} = P\{S_{k+1} \mid S_k\}}$$

### Markov Property

$$P\{S_3 = E \mid S_0 = A, S_1 = C, S_2 = D\} = P\{S_3 = E \mid S_2 = D\}$$

$$P\{S_3 = E \mid S_0 = A, S_1 = B, S_2 = D\} = \text{↑}$$

# Markov Decision Processes (MDPs)

Markov Processes evolve <u>autonomously</u>.
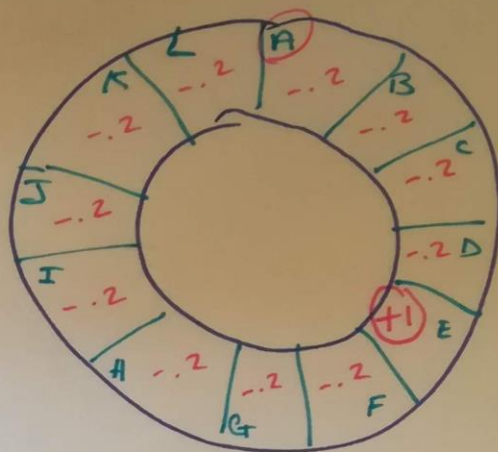Let's give Sisyphus a tiny bit of
free will:

  R: Throw rock counterclockwise
  L: Throw rock clockwise

Same throwing abilities for L & R

~~P{L}~~ P{$L_k$} unchanged



---

$\Rightarrow$ $T(A, L, B) = 0.25$  $T(A, L, C) = 0.5$  @$T(A, L, D) = 0.25$

   $\hookrightarrow$ $T(s, a, s') = P\{s' | s, a\}$          $T: S \times A \times S \rightarrow [0,1]$

          $\vdots$ $L = $action          set of possible actions.

$\Rightarrow T(A, R, L) = 0.25$  $T(A, R, K) = 0.5$  $T(A, R, J) = 0.25$

          $\vdots$

---

<u>Rewards</u>: $R: S \rightarrow \mathbb{R}$

For Sisyphus: $R(E) = +1$

          $R(s) = -0.2$ for $s \neq E$

Suppose Sisyphus
executes two actions

$a_1 = L$, $a_2 = L$
$\underset{k=1}{\llcorner}$       $\underset{k=2}{\llcorner}$

# Expectation

Suppose a random variable $X$ takes values from the set $\{c_1, c_2, \ldots, c_n\}$. The expected value of $X$ is defined as

$$E[X] = \sum_{i=1}^{n} P\{X = c_i\} \times c_i$$

---

Example: roll one die, $X$ shows

$$E[X] = \sum_{i=1}^{6} \frac{1}{6} \times i = 3.5$$

$\downarrow$

all values equally likely.

---

Intuition: Perform this many times

Average of $X$'s $\longrightarrow E[X]$

# Generalization

$$E[X_1 + X_2] = \sum_i \sum_j P\{X_1 = c_i, X_2 = c_j\}(c_i + c_j)$$

Two dice

$$E[X_1 + X_2] = \sum \sum \frac{1}{36}(i+j) = 7$$

---

$$E\left[\sum X_i\right] = \sum P\{X_i = c\} \sum(c_i)$$

---

For Sisyphus: $E[R(s_0) + R(s_1) + R(s_2)]$
given $a_1 = L$, $a_2 = L$ ??

We know

$\begin{cases} R(s_0) = -0.2 \text{ because } s_0 = A \\ R(s_1) = -0.2, \ s_1 \in \{B, C, D\} \\ R(s_2) = \begin{cases} +1 & s_2 = E \\ \\ -0.2 & \text{Else} \end{cases} \end{cases}$

# Expected return

Define return $r_h = \sum_{i=0}^{h} R(s_i)$. $\longleftarrow$ (✳)

$r_2$ for sisyphus $= R(s_0) + R(s_1) + R(s_2)$

$E[r_2] = E[R(s_0) + R(s_1) + R(s_2)]$

$\overset{P\{s_3 = E\}}{\underset{}{}}$ $\overset{P\{s_3 \neq E\}}{\underset{}{}}$

$= P\{r_h = .6\} \times .6 + P\{r_h = -.6\} \times -.6$ ✳

because $r_h$ has only two possible values

$-0.2 + -0.2 + -0.2 = -0.6 \longleftarrow s_2 \neq E$

$-0.2 + -0.2 + 1 = 0.6 \longleftarrow s_2 = E$

To arrive $s_2 = E$

| $s_0$ | $s_1$ | $s_2$ | Probability |
|-------|-------|-------|-------------|
| A | B | E | $(0.25) \times (0.25) = 0.0625$ |
| A | C | E | $(0.5) \times (0.5) = 0.25$ |
| A | D | E | $(0.25) \times (0.25) = 0.0625$ |

~~○○○○~~

$P(s_3 = E) = 0.375$

---

$P\{s_3 \neq E\} = 1 - P\{s_3 = E\}$

$= 0.625$

---

This is great for finite horizons (i.e. only consider a finite # of stages).

Suppose $h \to \infty$ ??

To deal with this, use discounted rewards.  $\underset{\text{factor}}{\overset{\text{discount}}{}}$

$r_h = \sum_{i=0}^{h} \gamma^i R(s_i)$

for $0 < \gamma < 1$.

---

$\lim_{h \to \infty} \sum_{i=0}^{h} \gamma^i R(s_i) \leq \sum_{i=0}^{\infty} \gamma^i R_{max}$

$\sum_{i=0}^{\infty} \gamma^i = \frac{1}{1-\gamma}$ $\quad 0 < \gamma < 1$

$\sum_{i=0}^{\infty} \gamma^i R(s_i) \leq \frac{R_{max}}{1-\gamma}$