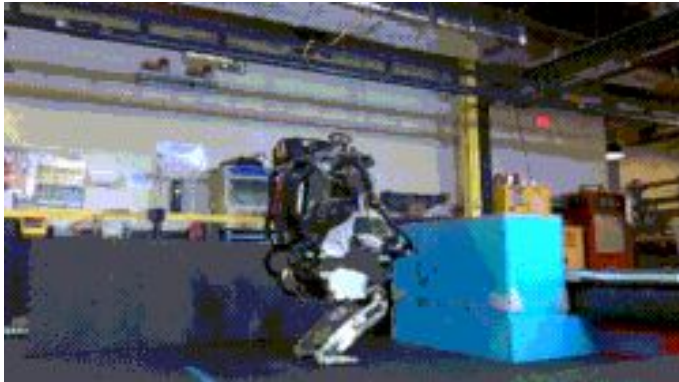


Crash Course in Reinforcement Learning

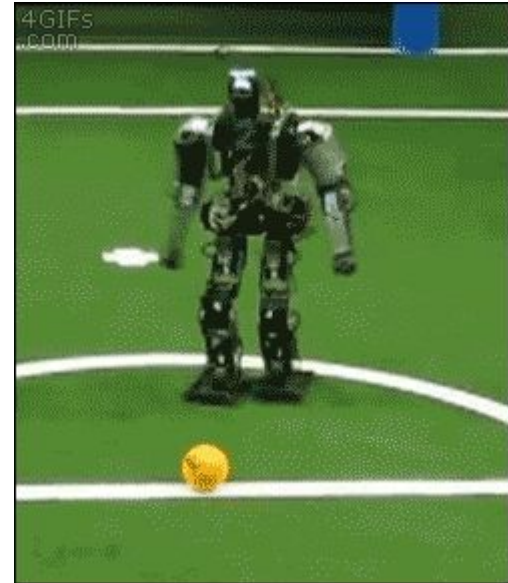
Wifi: welovecode

What is reinforcement learning?

- Teaching an agent to interact with an environment in order to achieve a specific outcome
- Why is it important?



@CLDSPN



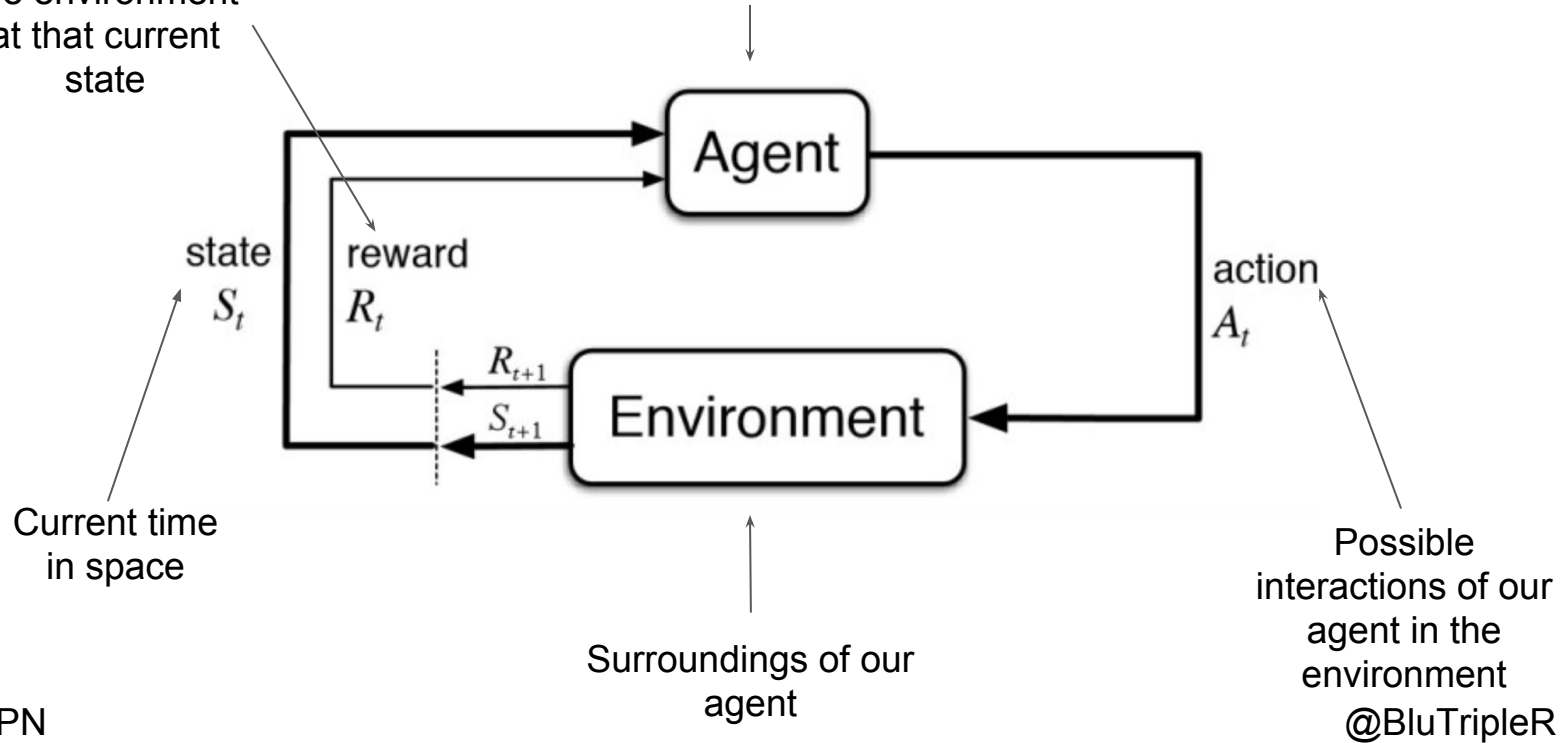
@BluTripleR

Why not just use supervised learning?

- Must build a dataset
- Learning from humans will only achieve the same level as humans
- Supervised algorithms usually require immediate feedback
 - Make a prediction and cross referencing with a label
- However, this is not the case with RL
 - Agent could be taking really useful actions, but then just fall at the last hurdle

Feedback to our agent based on actions taken in the environment at that current state

Entity being trained



@CLDSPN

@BluTripleR

State, Action, Reward

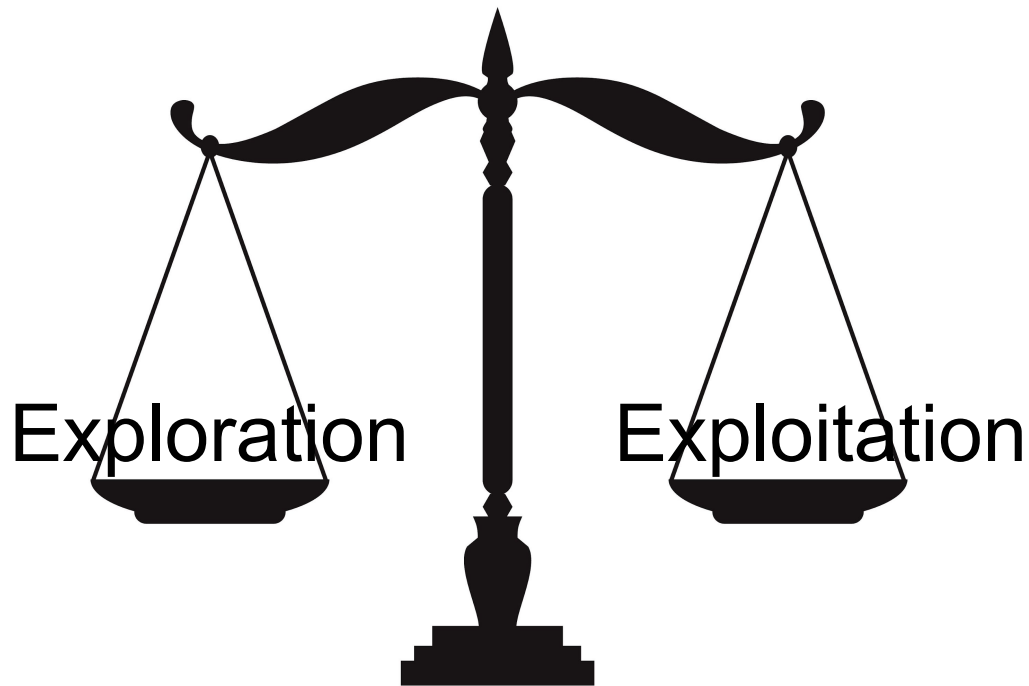
We want to teach our agent a desired behaviour.

We can hint at desired behaviour by ‘bread-crumbling’ rewards across states based on actions in the environment.

This is known as maximizing the cumulative reward

$$G_t = R_{t+1} + R_{t+2} + \dots$$

Exploration VS Exploitation



Exploration VS Exploitation

- Exploration- Curiosity
 - Discovering more information about the environment
- Exploitation- Insular
 - Top priority is seeking the highest reward

Quality Tables

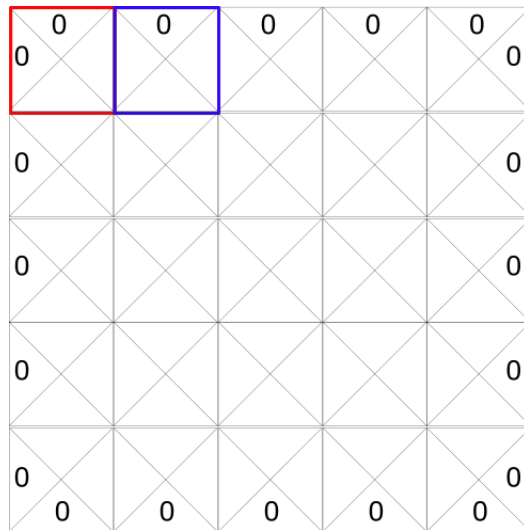


@CLDSPN

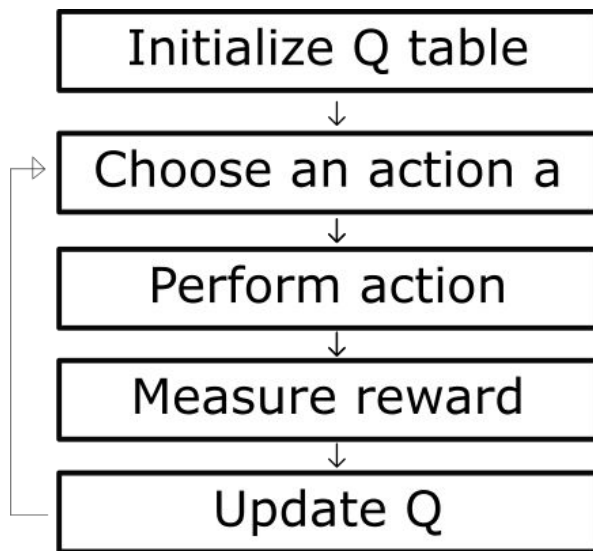
@BluTripleR

Q table

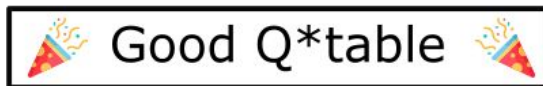
- Map environment landscape
- Every interaction gets logged for future reference
- Overtime agent understands the environment more and more



Q table update process

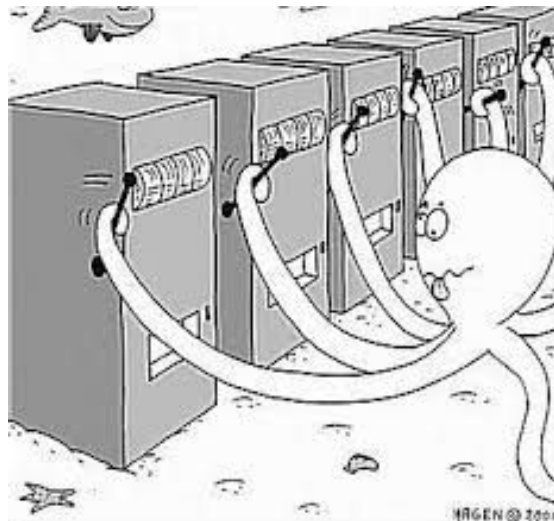


At the end of the training



$$\underbrace{NewQ(s, a)}_{\text{New Q value for that state and that action}} = \underbrace{Q(s, a)}_{\text{Current Q value}} + \underbrace{\alpha}_{\text{Learning Rate}} [\underbrace{R(s, a)}_{\text{Reward for taking that action at that state}} + \underbrace{\gamma}_{\text{Discount rate}} \underbrace{\max Q'(s', a') - Q(s, a)}_{\text{Maximum expected future reward given the new } s' \text{ and all possible actions at that new state}}]$$

Multi Armed Bandit



@CLDSPN

@BluTripleR

Multi Armed Bandit

- Find the action with the greatest amount of reward, while still earning a reward during this exploration phase
- Turbo charged A/B test
 - A/B testing- Used to gauge user preference
 - Problem- Separates exploration and exploitation

About Triptease

- SaaS startup building industry-leading software for the hotel industry
- Co-founded in 2015
- Offices in London, New York & Singapore
- \$24m raised with backing from BGF

TRIPTEASE

