

Reinforcement Learning

- MDP $\langle S, A, P, R, \gamma \rangle$ (POMDP)
- Bellman equations (V^π, Q^π, V^*, Q^*)
 - π^* from V^*/Q^*
- Dynamic Programming:
 - Policy evaluation, Policy iteration, Value iteration
 - update rule
 - 2 steps
 - update rule
 - Async DP, function approx, divergence
 - contraction proof (Bellman expectation & optimality backups)
- Model-free prediction:
 - MC & TD policy evaluation
 - TD(n) & TD(λ)
- Model-free control:
 - MC policy iteration, Sarsa, Sarsa(λ)
 - Q-learning, Deep Q-learning (exp replay & stationary targets)
- Policy gradient:
 - $\nabla_\theta J(\theta) = \dots$
 - Reinforce (baseline), Actor-Critic
- Model-based:
 - Sample-based planning; forward plan, simulation based plan, MCTS
 - Dyna
- Advantage actor-critic
- Exploration vs exploitation:
 - multi-armed bandit
 - minimize regret
 - UCB, confidence in face of uncertainty

RL

$$V^\pi(s) = \sum_a \pi(s, a) q^\pi(s, a)$$

$$= \sum_a \pi(s, a) \left[R_s^a + \gamma \sum_{s'} P_{ss'}^a V^\pi(s') \right]$$

$$q^\pi(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V^\pi(s')$$

$$= R_s^a + \gamma \sum_{s'} P_{ss'}^a \sum_{a'} \pi(s', a') q^\pi(s', a')$$

$$V^*(s) = \max_a q^*(s, a)$$

$$= \max_a R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s')$$

$$q^*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s')$$

$$= R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} q^*(s', a')$$

Policy evaluation:

• iterate Bellman exp eqn

Policy iteration:

• Policy evaluation + Greedy Policy improvement

Value iteration:

• iterate Bellman optimality eqn

Bellman
(expectation)
eqns

Bellman
Optimality
eqns

can do async,
or with priority
(by Bellman error)

can fit
into
policy
iteration too

MC Policy evaluation:

$$V^\pi(S_t) = V^\pi(S_t) + \alpha (G_t - V^\pi(S_t))$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T$$

TD Policy evaluation:

$$V(S_t) = V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

TD error, $G_t^{(n)} - V(S_t)$

equivalent can be
done w/ $Q(s, a)$

$$TD(n) G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$TD(\lambda) G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} = R_{t+1} + \gamma [(1-\lambda)V(S_{t+n}) + \lambda G_{t+1}^\lambda]$$

- Sarsa: policy iteration w/ TD policy evaluation

$$Q(s, a) = Q(s, a) + \alpha [R + \gamma Q(s', a') - Q(s, a)]$$

(Sarsa (λ) uses TD(λ) error, to do online need eligibility traces)

Expected Sarsa

- Q-Learning: off-line Sarsa

$$\rightarrow Q(s, a) = Q(s, a) + \alpha [R + \gamma \sum_a \pi(a|s') Q(s', a') - Q(s, a)]$$

for greedy π ,

$$Q(s, a) = Q(s, a) + \alpha [R + \gamma \max_a Q(s', a') - Q(s, a)]$$

- DQL:

experience replay stabilizes gradients

target network improves stability (stationary targets \Rightarrow SL problem)

- Policy Gradient:

$$\nabla_{\theta} J(\theta) = E [Q^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s)]$$

- REINFORCE:

Sample G_t from full episode (Monte-Carlo Policy Gradient)

$$\Theta_{R+1} = \Theta_R + \alpha \sum_t G_t \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$$

- Actor-Critic:

policy evaluation \rightarrow Critic: $v(s)$ or $q(s, a)$, update w/ TD (or TD(τ)/TD(λ))

policy gradient \rightarrow Actor: π_{θ} , update by policy gradient
(REINFORCE update, but G_t replaced with TD-error)

- Dyna:

every step training from reality, learn N steps from model
(experience buffer is a simple type of model)

- MCTS:

selection + expansion w/ tree policy, simulation w/ default policy,
update $Q(s, a)$ on path to root

- Regret:

$$L_t = \sum_{a \in A} v^* - q(a, t) = \sum_a N_t(a) \Delta_a \quad (\Delta_a = v^* - q(a, t))$$

for ϵ -greedy

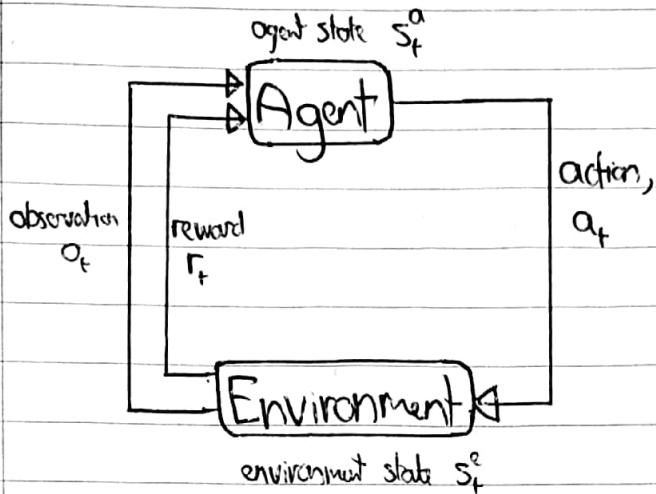
$$\rightarrow E[v^* - q(a_t, t)] \geq \frac{\epsilon}{|A|} \sum_a \Delta_a \Rightarrow \text{linear in } t \text{ w/ constant } \epsilon$$

- UCB:

$$P(E(X) > \bar{X}_t + u) \leq e^{-2tu^2}$$

$$P(q(a) > Q_t(a) + U_t(a)) = p \leq e^{-2N_t(a) + U_t(a)^2} \Rightarrow U_t(a) = \sqrt{\frac{-\log p}{2N_t(a)}}$$

Introduction to Reinforcement Learning



• Reward hypothesis:
any goal can be formalized as outcome of maximizing a cumulative reward

- o_t function of s_t^e
- s_t^a function of history; sequence of (o_t, a_t, r_t)
- a_t function of s_t^a
- r_t function of s_{t+1}^e and a_{t+1}

• Markov state:

$$S_t \text{ Markov} \Leftrightarrow P(S_{t+1} | S_t) = P(S_{t+1} | S_{1:t})$$

• future independent of past given present

• even if environment state is Markov, observation may not be (partially observable environment)

• Agent components:

$\pi(a|s)$ • Policy: map agent state to action

$q_\pi(s, a)$, $v_\pi(s)$, • Value function: expected future reward (discounted by γ)
depends on policy

$s', r = m(s, a)$ • Model: predict next state or reward

• can be used to help derive good policy

• Some combination of these 3 give a RL agent

Markov Decision Processes

• Markov process:

• tuple, $\langle S, P \rangle$

• $S = \text{finite set of states}$

• $P_{S,S'} = P(S_{t+1} = S' | S_t = S)$

• Markov reward process:

• tuple $\langle S, P, R, \gamma \rangle$

• $R_S = E(R_{t+1} | S_t = S)$

• $\gamma = \text{discount factor } (0 \leq \gamma \leq 1)$

• Return is discounted future reward

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

• Value is expected return

$$V(S) = E(G_t | S_t = S)$$

Bellman eqn

$$V(S) = E(R_{t+1} + \gamma V(S_{t+1}) | S_t = S) \leftarrow \text{for MRP}$$

$$V = R + \gamma P V$$

$$\Rightarrow V = (I - \gamma P)^{-1} R \quad (O(n^3) \text{ inversion; unfeasible for large state space})$$

• Markov Decision Process:

• tuple $\langle S, P, A, R, \gamma \rangle$

• A is finite set of actions

• $P_{S,S'}^a = P(S_{t+1} = S' | S_t = S, A_t = a)$

• $R_S^a = E(R_{t+1} | S_t = S, A_t = a)$

• Gully defining agent behaviour

• Policy is distribution over actions given states \leftarrow stationary

$$\Pi(S, a) = P(A_t = a | S_t = S)$$

• Markov Reward Process:

• State and reward sequence from following policy Π in an MDP

• Value function

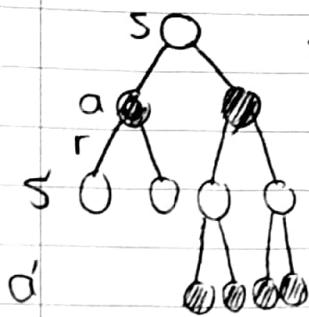
$$V^\Pi(S) = E_\Pi(G_t | S_t = S) \leftarrow \text{expected } G_t \text{ if follow policy } \Pi$$

$$q^\Pi(S, a) = E_\Pi(G_t | S_t = S, A_t = a) \leftarrow \text{take action } a, \text{ then follow policy } \Pi$$

• Bellman eqns:

Markov Decision Processes

- Bellman equations:



$$v^\pi(s) = \sum_{a \in A} \pi(s, a) \left[R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^\pi(s') \right]$$

$$q^\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(s', a') q^\pi(s', a')$$

$$v^\pi(s')$$

- v^* & q^* are max v^π & q^π over all policies
- solved MDP when found v^*/q^*
- optimal policy $\pi^* \geq \pi$ (achieving $v^* = v^*$, $q^* = q^*$)
 $(\pi \geq \pi \Leftrightarrow v^*(s) \geq v^\pi(s) \text{ for all } s)$
- $\pi^*(s, a) = \begin{cases} 1 & \text{if } a = \arg \max_a q^*(s, a) \\ 0 & \text{else} \end{cases}$

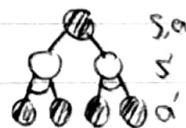
- Bellman Optimality eqns:

$$v^*(s) = \max_a q^*(s, a)$$

$$q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^*(s')$$

$$v^*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^*(s')$$

$$q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_a q^*(s', a)$$



- non-linear; no closed form solution in general; need to iterate
- Partially Observable MDPs

$$(S, A, O, P, R, Z, \gamma)$$

state set
of observations

observation function
 $Z_s^a = P(O_{t+1} = o | S_t = s, A_t = a)$

- Hidden Markov model with actions

Planning by Dynamic Programming

• Dynamic Programming:

- break complex problem into subproblems, solve those & combine solutions to solve complex problems

see Bellman eqns

value fn
store & reuse
solutions

- good when:
 - complex solution combination of simpler solutions
 - subproblems overlap/repeat
- \Rightarrow good for MDPs!

(Need full knowledge of MDP; then use DP to plan)

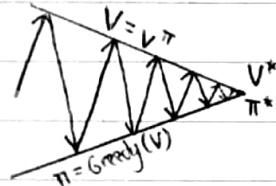
• Iterative policy evaluation:

- iterab $V_{k+1}(s) = \sum_{a \in A} \pi(s, a) \left[R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right]$

• Policy iteration:

- estimate V^π

- using iterative policy evaluation



- generate $\pi' \geq \pi$

- greedy policy on V^π

- improves as $V^\pi(s) \leq q^\pi(s, \pi'(s)) = V^{\pi'}(s)$

- if no change, $V^\pi(s) = \max_{a \in A} q^\pi(s, a) \Rightarrow$ Bellman optimality satisfied

• Value iteration:

- iterab $V_{k+1}(s) = \max_{a \in A} \left[R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right]$

- once V^* found, can derive $\pi^* = \text{Greedy}(V^*)$

- all $O(mn^2)$ per iteration (m actions, n states)

($O(m^2n^2)$ if based on q -function)

• Asynchronous DP:

- above have been done on all states per iteration; prohibitively expensive for large state spaces

- In-place DP:

- update as we go through states

- Prioritised sweeping:

- update (large) Bellman error = $\max_{a \in A} [R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s')] - V_k(s)$

- Real time DP:

- follow an agent in the MDP, updating states we visit

- Approximate DP:

- use a function approximator for value function, train to minimise square loss over a sampled mini-batch of (or every) states

- Divergence:

- Bootstrapping + General Function Approximation + not matching transition dynamics of MDP

Using state distribution

⇒ possibility to diverge (rare in practise)

- Contraction:

(& optimality backup)

- Bellman expectation backup is a contraction; brings value function closer by at least γ

- ⇒ by contraction mapping theorem all 3 algorithms converge to optimised policy or value function

improve value estimate by using estimate in future states

Model-Free Prediction

- Estimate value function of unknown MDP

- Monte-Carlo Policy Evaluation:

$$\cdot V^\pi(s) = E_\pi[G_t | S_t = s]$$

• Compute expectation by full rollouts following policy π

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$

learning rate from 1 full episode following π

- Temporal Difference:

• Bootstrap:

$$V(S_t) \leftarrow R_{t+1} + \gamma V(S_{t+1})$$

$$\text{TD error} \rightarrow \delta = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

• only perform 1 step, then update estimate using estimate

+ don't need to rollout to end

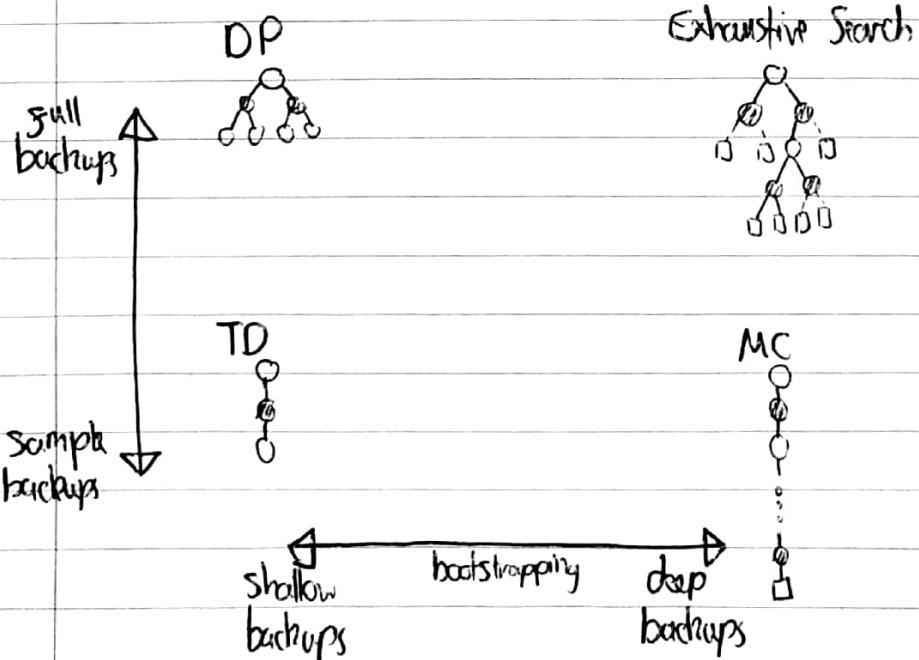
+ works in never ending environments

+ lower variance

- biased (a little)

- sensitive to initial value function

+ exploits Markov property; converges to max likelihood Markov model (Monte-Carlo method doesn't)



• TD(λ):

• use λ to decide 'degree' of bootstrapping

$$v(S_t) \leftarrow v(S_t) + \alpha (G_t^\lambda - v(S_t))$$

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}, \quad G_t^{(n)} = \sum_{i=0}^{n-1} \gamma^i R_{t+i+1} + \gamma^n v(S_{t+n})$$

$$G_t^\lambda = R_{t+1} + \gamma \left[(1-\lambda) v(S_{t+1}) + \lambda G_{t+1}^\lambda \right]$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
observe continu bootstrap continu weight
reward w/ weight w/ weight weight λ

• need full rollouts: can use eligibility traces to allow TD(λ) online

• $\lambda=1$ roughly equivalent to MC

• converges with tabular or linear value function approximation

$$(TD(n) \Rightarrow v(S_t) \leftarrow v(S_t) + \alpha (G_t^{(n)} - v(S_t)))$$

\uparrow
n-step
then bootstrap

Model-Free Control

- Optimize value function for unknown MDP

- Multi-armed bandit:

- A = actions space

- R_t is based on action $a_t \in A$

- how to maximize $\sum_t R_t$?

- need to trade-off exploration vs exploitation

- e.g. ϵ -greedy: prob $(1-\epsilon)$ take best-seeming action, prob ϵ take random action

- Policy iteration:

- estimate Policy evaluation:

- Monte-Carlo policy evaluation, $Q \approx Q^\pi$

- Policy improvement:

- ϵ -greedy policy using new Q

- Greedy In The Limit w/ Infinite Exploration: (GLIE)

- if all state-action pairs visited infinite times, policy converges on a greedy policy

- ϵ -greedy is GLIE if $\epsilon_k = 1/k$ (ϵ decays)

- \Rightarrow decay ϵ this way w/ policy iteration

- Sarsa:

- use TD learning instead of MC for policy evaluation

$$Q(s, a) \leftarrow Q(s, a) + \alpha \delta$$

$$\delta = R + \gamma Q(s', a') - Q(s, a)$$

change to $G_t^{(n)}$ for n-step backup

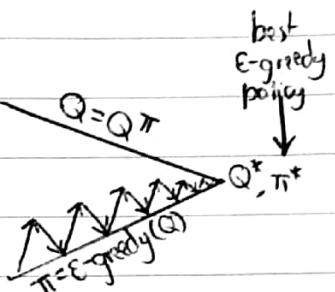
- Sarsa(λ):

updated for every s, a $\rightarrow Q(s, a) \leftarrow Q(s, a) + \alpha \delta_t E_t(s, a)$

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$

$$E_t(s, a) = 0$$

$$E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + I(S_t = s, A_t = a)$$



Off-policy
Learn about policy π by sampling policy μ

• Q-learning:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)]$$

• if $\mu = \pi$: Expected Sarsa

• if $\pi(S_{t+1}) = \arg \max_a Q(S_{t+1}, a)$ & μ ϵ -greedy wrt $Q(S, a)$
target simplifies to

$$R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$$

• converges to opt Q^* if take every action in every state infinite times

• Deep Q-Learning:

• To improve performance:

• Experience Replay:

• Store many previous (S, a, r, s') transitions, & update using minibatch to stabilise gradients

• Don't just do recent part, or we'll overfit to it

• Stationary targets:

• target depends on Q-network; changes every iteration

• instead, use target network, fixed for a while; problem closer to SL, improves stability

Policy Gradient

- Learn policy only; don't bother with value function or model
 - parameterise policy, $\pi_\theta(a|s) = P(a|s, \theta)$
 - + better convergence
 - + can learn stochastic policies
 - + policy could be simple whilst value/model complex
 - local optima & does not generalize

• Policy gradient:

- update parameters by gradient ascent

$$\Delta \theta = \alpha \nabla_\theta J(\theta), \quad J(\theta) \text{ is policy objective function}$$

eg avg reward per time step,
Average state value, or
initial state value

- need to compute

$$\nabla_\theta J(\theta) = \nabla_\theta E_d(V^{\pi_\theta}(S))$$

aka
contextual
bandit

- eg, 1 step case

$$\nabla_\theta J(\theta) = \nabla_\theta E(R(S, A))$$

can compute
this using
MC samples

$$\Rightarrow \theta_{t+1} = \theta_t + \alpha R_{t+1} \nabla_\theta \log \pi_{\theta_t}(A|S)$$

score function
trick

- also applies to long term reward, using $q^\pi(S, a)$

$$\nabla_\theta J(\theta) = E \left[\nabla_\theta q^{\pi_\theta}(S, A) \nabla_\theta \log \pi_\theta(A|S) \right] \quad (\text{expectation over } S \text{ & } A)$$

- don't need to know dynamics

- eg, trajectory T w/ return $G(T)$

$$\nabla_\theta J(\theta) = E \left[G(T) \nabla_\theta \sum_{t=0}^T \log \pi_\theta(A_t|S_t) \right]$$

• REINFORCE:

- Monte-Carlo Policy Gradient

learn from policy
evaluation

- for each episode:

$$\theta \leftarrow \theta + \alpha \sum_t G_t \nabla_\theta \log \pi_\theta(A_t|S_t)$$

similar idea to
target function

- problem: high variance

$$\text{use a baseline } b(S) = V_{\pi_\theta}(S) \approx V_\pi(S)$$

$$\theta \leftarrow \theta + \alpha \sum_t (G_t - b(S_t)) \nabla_\theta \log \pi_\theta(A_t|S_t)$$

- Actor-Critic:

- Critic:

- update η of V_η by TD

bootstrapping to
reduce variance

- Actor:

- update Θ by policy gradient

- for each step:

- Sample reward R_s^A, S'

- sample $A' \sim \pi_\Theta(S')$

- $\delta \leftarrow R_s^A + \gamma V_\eta(S') - V_\eta(S)$

TD error

- $\eta \leftarrow \eta + \beta \delta \nabla_\eta V_\eta(S)$

TD(0)

- $\Theta \leftarrow \Theta + \alpha \delta \nabla_\Theta \log \pi_\Theta(S, A)$

Policy Gradient

- $A \leftarrow A', S \leftarrow S'$

- can reduce bias by using TD(n) instead of TD(0)

$$\delta_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n V_\eta(S_{t+n}) - V_\eta(S_t)$$

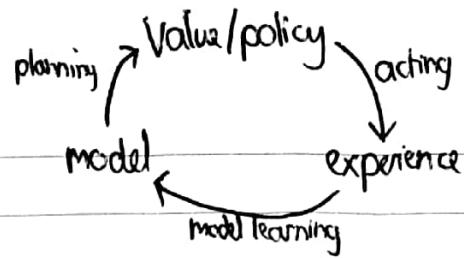
- Continuous actions:

- as only dealing with parameters on policy, easy generalise to continuous action space

- Gaussian policy common; $a \sim N(\mu(s), \sigma^2)$

$$\nabla_\Theta \log \pi_\Theta(S, a) = \frac{1}{\sigma^2} (a - \mu(s)) \nabla \mu(s)$$

Planning & Models



• Model:

- $M_n = \langle P_n, R_n \rangle$, approximates state transitions & rewards of an MDP $\langle S, A, P, R \rangle$

• Model learning:

- estimate model from experienced transitions $\langle S_t, A_t \rangle \rightarrow \langle R_{t+1}, S_{t+1} \rangle$
- SL problem; learn P & R functions
- table based, linear, neural network...

• Planning with model:

- solve MDP $\langle S, A, P_n, R_n \rangle$ using planning algorithms
- Value iteration, policy iteration...

Problem:
performance limited
to optimal policy for
MDP $\langle S, A, P_n, R_n \rangle$; not
optimal for true MDP

• Sample based planning:

- sample experience from model & apply model-free RL to this (MC, Sarsa, Q-learning)

• Model-based vs Model-free:

- can be fuzzy distinction; eg. experience replay buffer; effectively stores a model & samples from this to train
- can store experience in more complex forms too
- can do, eg. prioritised experience replay; train on transitions with largest Bellman error

• Dyna:

- learn using real & simulated experience; don't distinguish the two

• initialise $Q(s, a)$ & Model (s, a)

• do:

 sample a using $Q(s, a)$ (eg. ϵ -greedy)
 observe R_s^a & s'

 update $Q(s, a)$ using Q-learning

 update Model $(s, a) \leftarrow (s', R_s^a)$

 repeat N times:

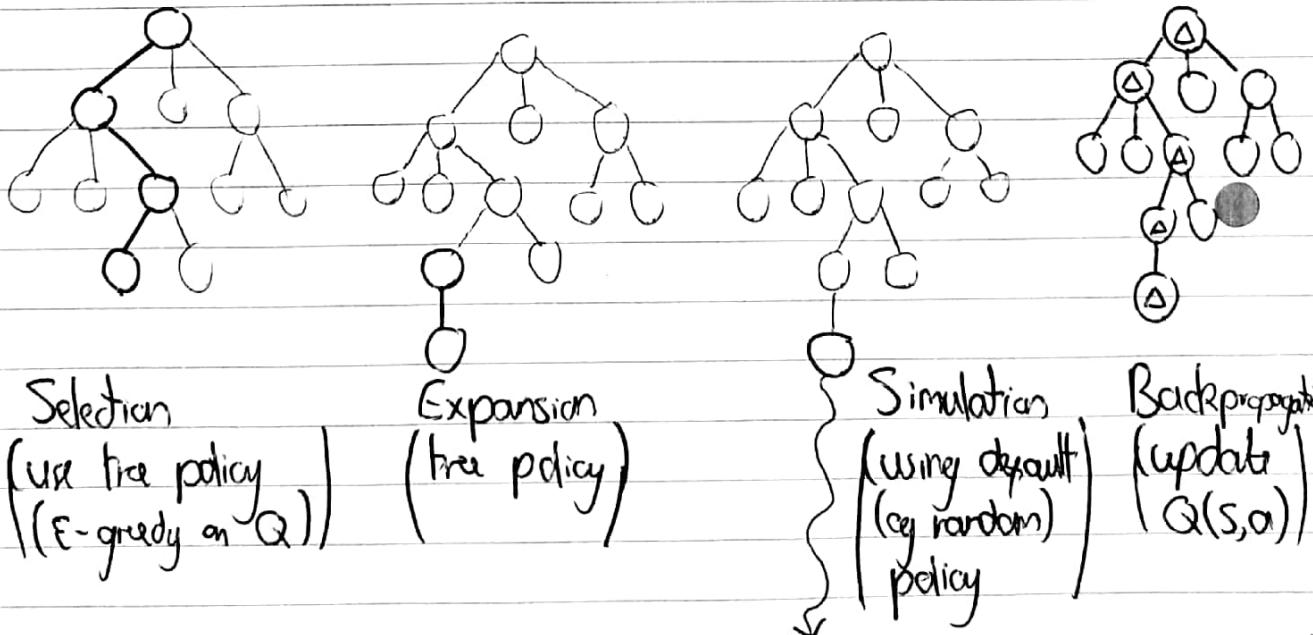
 randomly choose s , pick a chosen in s last time encountered,
 find s', R_s^a using Model, update $Q(s, a)$ using Q-learning

• Planning:

- given fixed model, how to use optimally to select next action
- forward search:
 - search from now to termination to determine best action
 - use model
- simulation-based search:
 - simulate episode from now on
 - use model-free RL on the episodes
 - MC control \rightarrow MC search
 - Sarsa \rightarrow TD search

• Monte-Carlo Tree Search:

- simulate K episodes, using current simulation policy π
- build search tree from this
- evaluate $Q(s, a)$ (for all s, a ^{leavy} ~~visited~~) by mean return from episodes from s, a
- select action by $\pi \text{ argmax}_a Q(s_t, a)$



Repeat K times

Then select action

Building Agents

$$A_t = \text{agent}(R_t, O_t)$$

(policy is a simple agent, only looking at O_t)
(learning done through looking at R_t too)

• e.g. neural Q-learning:

• network $q_\theta: Q_t \rightarrow \mathbb{R}(q[a_1], \dots, q[a_m])$

• ϵ -greedy policy: $\pi(q_t) \rightarrow a_t$

• Q-learning: squared loss

• optimizer to minimize loss (SGD, RMSProp, Adam)

• Advantage Actor Critic:

• Representation: $(S_{t+1}, O_t) \rightarrow S_t$

• network $V_\theta: S \rightarrow V$

• network $\pi_\theta: S \rightarrow \pi$

• copy π^m of π_θ as policy

$S_t^m \rightarrow A_t^m$

n-step TD loss on V_θ

n-step REINFORCE loss on π_θ

optimizes for each

• Exploration vs Exploitation:

• aim to minimize total regret

$$L_t = \sum_{i=1}^t V^* - q(a_i)$$

$$V^* = \max_{a \in A} q(a)$$

regret =
opportunity cost
for 1 step

$$q(a) = E[R_t | A_t = a]$$

• ϵ -greedy: prob $(1-\epsilon)$: greedy action, prob ϵ : random action

$$E(V^* - Q_t(a_t)) \geq \frac{1}{|A|} \sum_{a \in A} \Delta_a \quad (\Delta_a = V^* - q(a))$$

hard to
find decay
rate

\Rightarrow constant $\epsilon =$ linear total regret

\Rightarrow decaying $\epsilon =$ logarithmic asymptotic total regret

• Optimism in the face of uncertainty:

◦ when very uncertain about reward from action, better to explore that one

• Upper confidence bounds:

- large $U_t(a)$ → \Rightarrow uncertain about a
- estimate upper confidence bounds $U_t(a)$ for each action (eg 95% confidence bound) $(P[q(a) > Q_t(a) + U_t(a)]) \approx$
 - select action with highest UCB
- \Rightarrow Optimistic in face of uncertainty

• UCB for Multi-Armed Bandit:

$$P[E(X) > \bar{X}_t + u] \leq \exp(-2t u^2)$$

$$\text{for iid } x_1, \dots, x_t, \bar{X}_t = \frac{1}{t} \sum x_i$$

$$\Rightarrow P[q(a) > Q_t(a) + U_t(a)] \leq \exp(-2N_t(a)U_t(a)^2) = p$$

$$U_t(a) = \sqrt{\frac{-\log p}{2N_t(a)}}$$

decay p over time, e.g. $p = \frac{1}{t}$
 $(\Rightarrow$ select optimal action as $t \rightarrow \infty$)

$$\Rightarrow a_t = \arg \max_{a \in A} Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}} \quad (c > 0)$$

(UCB algorithm achieves logarithmic total regret)