

1. Introduction & Lecture 1
How do we build systems that solve tasks for which humans need intelligence?

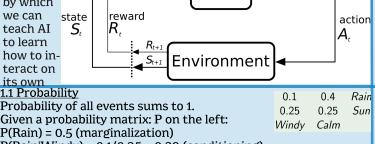
ML: The answer to the AI question: Methods, algs, data structures that learn to solve such tasks leveraging data.

Big Data: Encompasses methods to handle and understand large data sets - composed of Data Science (how to analyse the world a data driven way) and Data Engineering (how to collect, store, process, clean and maintain large data).

AI usually refers to an intelligent system that embodies an entire practical system (self driving car) while ML is focused on the algos/ learning mechanisms for building a system to solve a problem rather than engineering a system that solves problems.

Reinforcement Learning: Solves problems of control - choosing the best / optimal action at the right time.

The core of RL: Agent interacts with environment and gains knowledge, exploring and receiving rewards. Actions change the state of the env, we choose actions for max long-term reward.



1.1 Probability
Probability all events sums to 1.
Given a probability matrix: P on the left:
P(Rain)=0.5 (marginalization)

P(Rain|Wind)=0.1/0.35=0.29 (conditioning)
Should I bring an umbrella if its windy (cost=1) - getting caught

the rain with an umbrella is cost=2 (central)

Cost of umbrella = 1. Cost of no umbrella = 2*0.29 = 0.58

Bayes Theorem: $p(A|B) = p(B|A)P(A)/P(B)$

Boolean Logic: If A is true B is true. B is true - doesn't tell us abt

Bayes Logic: If A is true, B is true. B is true, so A is more likely

Plausibilities: (encodes prior info) 3 criteria apply to them:

1) The degrees of plausibility are represented by real numbers

2) These numbers must be based on the rules of common sense

3) Plausibilities are self-consistent when found via different methods, they take into account all available data, and discovering the plausibilities is reproducible.

Posterior: $p(\text{Parameters}) = p(\text{Data}|\text{Parameters}) * p(\text{Parameters}) / p(\text{Data})$
Likelihood / Evidence: $p(\text{Data}|\text{Parameters})$
Prior: $p(\text{Parameters})$ - our prior belief over possible param values

We can visualise structure of dependencies between random variables with directed edges of causation events (with edges having Ptarget/source) from source to target nodes and nodes with pure source nodes reflecting prior probabilities P[source].

Benefits of Bayes Logic: we can make inferences based on uncertain information. Probabilities can represent degrees of belief. We can make principled decisions - Is P(A)? Cost(A) v.s. A

Bayesian Decision Theory: $a^* \in \arg \max_a E[r(a)] = \arg \max_a \sum_{j=1}^J p(j, a)p(j)$

Making optimal decisions by maximizing a^* , a represents a decision/action, s is the state

We do this in RL, Optimal Control (Robots).

RL involves learning optimal control of an "a priori" unknown system with unknown rewards - only experience teaches us

1.2 Markov Processes

A Markov Process is a tuple (S, P, τ) . S is a set of states and $P(s')$ is a state transition probability matrix $P(s'|s) = \tau(s'|s)$.

Markov Property: A state is Markov iff $P(s'|s) = P(s'|s_{-1}, s_{-2}, \dots, s_0)$ - that is, the next state is only determined by the current state.

Square states are terminal, round are transient (can start there).

That's why $P(s'|s)$ is 1 (of course, as all transition probabilities sum to 1)

Stationary Markov Chains: If the $P(s'|s)$ don't depend on t, (and so only on origin/dest state) the Markov chain is stationary.

Markov Reward Process: A Markov chain that emits rewards.

Formed of a tuple (S, P, R) . P and S are as before.

R = $E[r|s] = E[r|s, t]$ is an expected immediate reward collected at time step t up to departing state s.

Expected total discounted reward is $\sum_{t=0}^{\infty} \gamma^t r_t$ where $\gamma \in [0, 1]$

The γ factor is how we discount the current value of future reward, v close to 0 leads to myopic evaluation (as looking at the return we ONLY care about the immediate reward - what's in the current state), v close to 1 is farsighted evaluation - we factor lots of distant rewards.

Most Markov Chains are discounted, i.e. $v < 1$:

• It avoids infinite returns in cyclic processes

• Helps take into account uncertainty about the future

• Immediate rewards can help us gain future ones (e.g. finance)

• Human & Animal decision making favours immediate rewards

• Sometimes $v = 1$ is all sequences terminate, or when they're equal length

State Value Function $v(s)$: The expected return R starting from state s at time t: $V = E[R|T=s]$ by definition (**expectation**)

This equals: $E[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | T=s]$

E [r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | T=s] = E[r_1 + \gamma V(s)| T=s]

Which represents two terms - the immediate reward r_1 and the discounted reward of the successor state $v(s|t)$.

Sum notation: $v(s) = R_s + \gamma \mathbb{E}[P(s')v'(s')]$ (n eqs, 1 for each state)

The vector equation is directly solvable:

$v = (1-\gamma)R + \gamma M$. But matrix inversion is expensive $O(n^3)$: direct

implementations are only good for small MRPs. Dynamic Programming, Monte-Carlo Evaluation, Temporal-Difference learning faster.

Policy: A policy $\pi(a|s) = P(A=a|S=s)$ is a conditional probability distribution (not deterministic) to execute an action a at given that one is in state s at time t.

The policy is only deterministic if for a given state only 1 a is possible: $\Pi(s) = 1$ and $\Pi(s) = 0$, $\forall a \in A$. Shorthand: $\Pi(s) = a$

1.3 Markov Decision Processes - like MRP but with an action

S: State Space, A: Action Space, $v(s) = \mathbb{E}[r|s]$ is discount factor, $\pi(a|s) = \Pi(a|s)$ (probability)

Deterministic: $a = \Pi(s)$ (indicator function)

1. State Value Function for MDPs: (A)

$$V'(s) = E[R_t | T=s] = E[\sum_{k=t}^{\infty} \gamma^k r_{t+k+1} | T=s]$$

Where R is a discounted total return and r_{t+k+1} is the reward

Backup Diagrams: Traces paths from state s to successor states s' given we chose action a. We can transfer state value info back to a state s' from successor states s' - for a back operation to find the right Value Func:

1) Find the probability of taking action a (governed by our policy):

2) The probability of transition to s' (transition probability: $P(s'|s, a)$)

3) The reward of each (s, a, s'): $R_{ss'}^a$

4) The value of s' (s') recursive definition:

$$V_{ss'}(s) = \frac{1}{1-\gamma} \mathbb{E}[r_{t+1} + \gamma V_{ss'}(s') | T=s]$$

5) $V(s) = \sum_{a} \sum_{s'} \pi(a|s) V_{ss'}(s)$

6) $V(s) = \sum_a \sum_{s'} \pi(a|s) P(s'|s, a) R_{ss'}^a + \gamma V(s')$

7) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s')]$

8) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma \mathbb{E}[V(s')]]$

9) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

10) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

11) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

12) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

13) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

14) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

15) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

16) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

17) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

18) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

19) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

20) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

21) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

22) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

23) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

24) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

25) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

26) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

27) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

28) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

29) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

30) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

31) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

32) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

33) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

34) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

35) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

36) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

37) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

38) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

39) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

40) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

41) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

42) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

43) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

44) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

45) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

46) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

47) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

48) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

49) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

50) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

51) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

52) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

53) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

54) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

55) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

56) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

57) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

58) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

59) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

60) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

61) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

62) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

63) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

64) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

65) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

66) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

67) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

68) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

69) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

70) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

71) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

72) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

73) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

74) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

75) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

76) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

77) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

78) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

79) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

80) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

81) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

82) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

83) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

84) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

85) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

86) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

87) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

88) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

89) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

90) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

91) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

92) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

93) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

94) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

95) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

96) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

97) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

98) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

99) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

100) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

101) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

102) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

103) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

104) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

105) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

106) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

107) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

108) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

109) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

110) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

111) $V(s) = \sum_a \pi(a|s) [R(s, a) + \gamma V(s)]$

Initialise replay memory D to capacity N
 Initialise action-value function Q with random weights
 for episode = 1, M do

 Initialise state s_t

 for $t = 1, T$ do

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(s_t, a)$

 Execute action a_t and observe reward r_t and state s_{t+1}

 Store transition (s_t, a_t, r_t, s_{t+1}) in D

 Set $s_{t+1} = s_t$

 Sample random minibatch of transitions (s_t, a_t, r_t, s_{t+1}) from D

 Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(s_t, a_j; \theta))^2$

 end for

end for

Optimization 2: Target Network

• Whenever we do an update step on a Q-network's state we also update "nearby" states (monolithic & generalisation)

• Problem: unstable training resulting from bootstrapping a continuous state space representation:

• We may have an unstable learning process because

 changes to $Q(s, a)$ change $Q(s', a')$ (for $s-s' < \delta$), which changes $Q(s', a')$ on the next update and back and forth.

Causes resonance where one update mirrors many times

• This can result in run-away bias due to bootstrapping, and bias may dominate the system. Q values then diverge.

• Involves TD error in target function changing often

Slowing things down (Responsible Damping) solves the often

• Have a separate target NN Q' , updated less often (e.g., every 1000 steps) with a copy of the latest learned weight parameters, controls this stability (relaxation time).

• Init two Q-nets, a main Q net and a target net Q'

• When calculating TD error use Q not Q'

• Infrequently set $Q' = Q$

• Gives the fluctuating Q relaxation time to settle before updating Q . Relax time can be hyperparam optimized.

Optimization 3: Clipping of Rewards:

• Different tasks have different reward scales, to have an RL model generalise across all of them we need to scale rewards

• We clip rewards: positives are set to 1, negatives set to -1

Optimization 4: Skipping Frames:

• Computers can simulate games (60Hz). Since games are designed for humans that react much more slowly, we don't need to react at 60Hz - games designed that way!

• Frame Skipping - for ATARI DQN we use only every 4 frames (so operate at 15Hz) and frame stacking uses past 4 frames as inputs) reducing computational costs and speeding up training time. Games run closer to human reaction pace.

6.1) Advanced DNNs - avoiding falling Q-Values while Training

High variance MDPs can confuse the Q learner. Sometimes a path gives really high reward, otherwise it's negative.

Since we update with a max over all actions of succ state:

$Q(s_t, a) \leftarrow Q(s_t, a) + \alpha_{t+1} + \gamma \max Q(s_{t+1}, a) - Q(s_t, a)$

we can have maximization bias: when taking the max over all actions with finite state we may overestimate.

Situation: the target network does 2 tasks:

1. We use the target network to identify the action with best Q

2. We also use it to find the Q value of this action

Problem of this: max Q value may be overestimated (var-bias problem) because an unusually high value from the main net Q doesn't mean there's an unusually high val from target net Q'

Solution: Double Q Learning

• We use a target network Q' for 1 but regular Q for 2 (or vice-versa). So the selection of the action with highest Q and the Q value are both independent. So we're less susceptible to the variance of each estimate. Reduces freq of overestimating the max Q value as it's less likely both nets overestimate the same action.

• DQN is also more realistic to the final values.

• Policy Gradient - find optimal policy without O/V function Examples: Policy-Based Algorithms (like Policy Gradient), REINFORCE, DDPG, Soft-Actor-Critic

Before we used value based methods. Here we instead use policy based methods which have a faster convergence and considered better for continuous / stochastic environments:

• We learn directly a parameterised policy π selecting actions without a value function with policy weight parameters θ

• We optimize by looking at traces that a policy $\pi(\cdot | \theta)$ rolls out to correlate them with the rewards they incur

• Probability to observe t depends on params of policy

$p(\tau) = p(s_1, a_1, \dots, s_T, a_T)$

The optimal

policy isn't from an optimal VF V^*

or Q^* but optimal

parameters θ^* :

aka those that give us a maximum average return:

$\theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)}$

Rationale: Continuous $\theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)}$

Parameterisation will be good

$\theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)}$

for continuous spaces which are otherwise too expensive for point-based approaches.

• The Policy Gradient Method:

• Uses gradient ascent to move towards the direction suggested by the gradient to find the best policy

• VFs can be used to learn policy weights but not for selection

• Policy Gradient Methods are methods for learning the policy weights using the gradient of a perf measure w.r.t weights

• Search directly for π^* for optimising params θ

infinite horizon case:

finite horizon:

$\theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)}$

Any supervised learning model can be used to learn θ

Finding the Gradient of a Policy:

• Notation: $J(\theta|s_t)$ is the probability the action a is taken in state s given policy weight vector θ

• The performance measure $J(\theta)$ is the value of the start state under the parameterised policy: $J(\theta) = V^*(s_0) = \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)}$

To optimise θ we need the deriv:

• derivatives on action selection (determined by $J(\theta|s, a)$)

• derivatives on the distribution of states $p(s)$, indirectly determined by $J(\theta|s, a)$

We can interpret the term $f(s, \theta) \cdot \partial_t$ as a curious form of supervised learning - where we fit the network output computed on the states $f(s)$ to the actions a .

It is curious because pure maximum-likelihood-type supervised learning (fitting states and actions to each other) is actually reward-weighted fitting in Policy-Gradient-based RL.

1. Finding the derivative: Finite Difference Policy Gradient

Vector diff'n of an M dim vector x is the partial deriv:

• Compute the partial deriv of the perf measure $\nabla_{\theta} f = J(\theta)$ w.r.t. θ using finite diff gradient approximation

• Then compute $J(\theta)$

Steps:

1. Loop over all M dims of θ with index k , set u_k as unit vector in the k th component, 0 elsewhere; e.g. $u_0 = [0 \ 0 \ 0 \ \dots \ 0]^T$

2. Compute partial ∂_t by small ϵ in k th dim only and compute $J(u_k)$

3. Then compute $J(0)$

These

• These can be interpreted as weighted avg over policies:

• Increase the weight of trajectories that are good in reward

• Decrease the weight of trajectories that are bad in rewards

2. Finding the Derivative: Direct Policy Gradients

We have the policy gradient theorem: $\theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)}$

Average return is approximated by empirical mean over T traces

$J(\theta) = \mathbb{E} \left[\sum_{t=1}^T r(s_t, a_t) \right]_{\pi_\theta(s_t, a_t) \sim p(\tau)} \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \quad (10)$

$\theta^* = \arg \max_{\theta} \mathbb{E}_{\pi_\theta} \left[\sum_{t=1}^T r(s_t, a_t) \right]$

• A convenient identity

$\nabla_{\theta} r(s, a) \nabla_{\theta} \log \pi_\theta(r) = \pi_\theta(r) \frac{\nabla_{\theta} \pi_\theta(r)}{\pi_\theta(r)}$

$J(\theta) = E_{\pi_\theta(r|\tau)}[r|\tau] = \int \pi_\theta(r|\tau) r d\tau$

$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} \pi_\theta(r|\tau) r d\tau = \int \pi_\theta(r|\tau) \nabla_{\theta} \log \pi_\theta(r) r d\tau = E_{\pi_\theta(r|\tau)}[\nabla_{\theta} \log \pi_\theta(r)|\tau]$

$\theta^* = \arg \max_{\theta} \mathbb{E}_{\theta} [\theta] = \log \text{both sides}$

$\frac{\partial}{\partial \theta} \log \pi_\theta(r|\tau) = \frac{\partial}{\partial \theta} \log \pi_\theta(s_{i,t}, a_{i,t}) = p(s_i) \prod_{k=1}^t \pi_\theta(a_k|s_{k+1}, a_k)$

$\nabla_{\theta} \log \pi_\theta(r|\tau) = \sum_{t=1}^T \log \pi_\theta(a_t|s_t) + \log p(s_i)$

$\nabla_{\theta} J(\theta) = E_{\pi_\theta(r|\tau)}[\nabla_{\theta} \log \pi_\theta(r|\tau)]$

We arrived at this final eq. a version of the policy gradient theorem, by expanding the state sequence from a fixed starting point s_1 . Generally we can obtain results for the stationary distribution of starting states.

REINFORCE: implements the Policy Gradient Theorem result:

• loop: improve the policy \rightarrow generate samples (run policy) \rightarrow model a fit to sample the return \rightarrow improve the policy...

1. sample $\{s^t\}$ from $\pi_\theta(s|s_t)$ (run the policy)

2. $\nabla_{\theta} J(\theta) \approx \sum_t (\sum_{s^t} \nabla_{\theta} \log \pi_\theta(a_t|s_t)) (\sum_t r(s_t, a_t))$

• Before we had to learn a VF through function approxim.

then derive a corresponding policy

But learning the VF could be intractable.

REINFORCE allows us to directly optimize policies to avoid this.

• REINFORCE suffers from high variance in samples

trajectories, stabilising model params is difficult

• Any erratic trajectory can cause a sub-optimal shift in policy distribution if it by chance gets low/high rewards

• To solve this we use a baseline or actor-critic.

• REINFORCE: first attempt at policy-grad with θ bias

but high variance is a problem. Future algorithms reduce variance by being smarter with correlating rewards with trajectories. Involve subtracting a baseline from reward (to reduce rewards), or using an advantage term

• Converges in the limit of large amounts of data

Causality: A way to reduce variance:

• Apply Causality to lower amt of total data (and so variability):

$\nabla_{\theta} J(\theta) = \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_\theta(a_t|s_t) (\sum_{t=1}^T r(s'_t, a'_t))$

• The policy at time $t > t'$ obviously can't affect the reward at t so rewrite the terms entering the sums as:

$\nabla_{\theta} J(\theta) = \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_\theta(a_t|s_t) (\sum_{t=1}^{t'-1} r(s'_t, a'_t))$

• We now have $J(\theta)$ as a function of the policy π_θ and not the policy weight θ .

• We can now update θ to $\theta + \alpha \nabla_{\theta} J(\theta)$

end for

2. DDPG Optimizations to make it work:

1. Exploration in continuous action spaces:

• Use Gaussian Distributed Noise that perturbs a mean action $\mu'(s) = \mu(s) + N$ - similar to discrete action e-greedy

2. Replay Buffers, Minibatches & Target Networks

controls the learning variability

3. Error-based learning on the mean-squared Bellman error

• MSBE (L) tells us how close we came to satisfying the BE:

$L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta_Q))^2$

4. Policy Gradient Based Learning on an Actor Network:

• We use a deterministic action version of the PG Theorem

• DDPG updates the target net in jumps, with many

remaining episodes being frozen. Doesn't work in continuous control; small policy changes must cause

small control changes.

• So instead, DDPG uses smooth updates on the param vec

that is shared by both the actor and the critic which

updated slowly but steadily over steps t:

$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

• The actor updates the policy in the air suggested by the critic (such as w/ policy gradients or updating μ).

• The critic estimates the value function. This could be the state-action value (Q function) or state-value (V function), (such as Function Approximation/DQN-style updates of Q).

A simple Q-driven policy-gradient actor-critic.

Algorithm 1 Q Actor Critic:

Initialize parameters s, θ, w and learning rates α_θ, α_w ; sample $a \sim \pi_\theta(a|s)$ for $t = 1, \dots, T$ do

 Sample reward $r_t = R(s_t, a)$ and next state $s'_t \sim P(s'|s_t, a)$

 Then sample the next action $a' \sim \pi_\theta(a'|s'_t)$

 Update the policy parameters: $\theta \leftarrow \theta + \alpha_\theta Q(s_t, a_t | \theta_Q)$

 Compute the correction (TD error) for action a at time t: $\epsilon_t = r_t + \gamma Q(s'_{t+1}, a'_{t+1} | \theta_Q) - Q(s_t, a_t | \theta_Q)$

 Update the action parameters of Q function:

$w \leftarrow w + \alpha_w \nabla_w Q(s_t, a_t | \theta_Q)$

Move to $a \leftarrow a' \text{ and } s \leftarrow s'$

end for

Reducing Variance: Baselines

• We are summing up our PG terms for 3 trajectories we evaluated: $\text{VLog|I|theta(t)}: 1/4, 1/5, 1/6$ and 0.4

• Our 3 trajectories have returns $r(t) = 100, 101, 102$ and so the variability is $\text{SD}(2, 102) = 10, 10.1, 10.2$

• If we want to subtract a baseline of $b = 100$ from the return $r(t)$ so $r(t) - b = 0$ \rightarrow 1 and 2, we obtain terms $1/4, 1/5, 1/6$ and 0.4

• We thus reduce variability by subtracting a good baseline

• Putting this into practice gives Advantage-Actor-Critic

• The key element is the Advantage Element. It tells us what the advantage amount is by picking a state in s. e.

• Q functions are decomposed into 2 pieces: the original state value function ($V(s)$) and the advantage value $A(s, a)$ of taking a specific action given a state: $Q(s, a) = V(s) + A(s, a)$

• Advantage func says how much better an action is compared to the others in s, while as we know the VF captures how good it is to be in this state.

$A(s, a) = Q(s, a) - V(s)$

• It's units are bits, e.g. a fair coin has the entropy of 1bit.

• Max entropy is when all outcomes xi are equally probable, min (zero) entropy is when outcome is deterministic.

For Continuous Random Variables:

$H(x) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i)$

• So the eval of an action is based on how good the action is - how much better it can be than the others

• Reduces the high var of policy net stabilizes training

• Reduces var without introducing bias in Policy Gradient

• Asynchronous Advance Actor Critic has multiple independent agents (nets) with own weights, which interact with a diff copy of the env in parallel. So can explore a bigger part of the action-space in much less time

• Workers train in parallel and update periodically a global network, which holds shared parameters.

• The updates are not happening simultaneously.

• After each update, agents reset parameters to those global net's & continue their independent exploration and training for n steps till they update again.

• Info flows from the agents to the global net but also between agents. As an agent resets weights by the global net, it interacts w/ the info of all other agents.

• Useful with large scale simulators, or multiple identical robotic facilities, to run algorithms in parallel.

Soft-Actor-Critic continued:

• Soft Actor-Critic has two Q-Value Functions, Q_{soft} and Q_{W} . This mitigates maximization bias in policy improvement.