

Elements common to all control task

s_t - state, a_t - Action r_t - Reward, Agent

Environment,

Markov decision process MDP

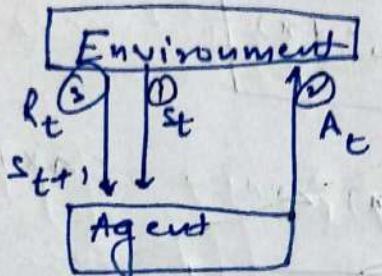
Discrete time

Time moves toward infinite interval
 $t \in [1, 2, 3, \dots]$

stochastic control process

future states depends only partially on action taken

It is based on decision making to reach the target state



set of possible states (s, a, r, p) set of actions that can be taken in each of the states

Probability of passing from one state to another for when taking each possible action

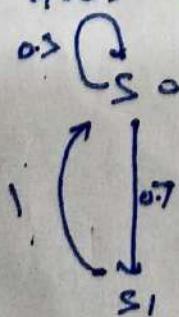
set of Rewards for each (s, a) pair

- The process has no memory

$$P(s_{t+1} | s_t = s_t) = P[s_{t+1} | s_t, s_{t-1}, \dots, s_0 = s_0]$$

- The next state depends only on the current state not on the previous state

- If a process meets this property, it is known as markov process



Different Action reward

Types of Markov decision process

Infinite finite
(game)

Infinite
(car driving)

Terminating P.M. ending condition like done

continuing

Trajectory and episode \rightarrow start to target

Trajectory :- elements that are generated when the agent moves from one state to another

$$\tau = s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots, r_T, s_T$$

Episode Trajectory from the initial state of the task to a terminal state

$$\tau = s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots, r_T, s_T$$

Reward vs Return

The goals of the task are represented by the rewards (R)

We want to maximize the sum of rewards

A short term reward can worsen long-term result

Return

Reward

r_t

$$G_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots + r_T$$

Discount factor (γ) The agent has no incentive to go to the goal through the shortest route

$$G_d = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^{T-t} r_T$$

We will multiply future rewards by a discount factor γ

$$G_d = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^{T-t} r_T$$

$$\gamma \in [0, 1]$$

higher the exponent more value of reward is today

If $\gamma = 0$, all future rewards will be 0

If $\gamma = 1$, rewards will not be discounted

- we want to maximize the long term discounted sum

Policy $\pi(s)$ function that decides what action to take in a particular state

$$\pi : s \mapsto A$$

Probability of taking action
a in state s

$$\pi(a|s)$$

stochastic

$$\pi(s) = [p(a_1), p(a_2), \dots, p(a_n)]$$

$$\pi(s) = [0.3, 0.2, 0.5]$$

Deterministic

$$\pi(s) \rightarrow a$$

$$\pi(s) = a$$

we want to maximize the sum of discounted rewards

we must find the optimal policy π^*

state values and action values

$$\text{state value } V_{\pi}(s) = \mathbb{E}[G_t | s_t = s]$$

$$V_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T | s_t = s]$$

following policy π

q-value of a state-action pair

$$q_{\pi}(s, a) = \mathbb{E}[G_t | s_t = s, A_t = a]$$

$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T | s_t = s, A_t = a]$$

following policy π

Bellman equation $v(s)$:- searching optimal policy

$$V_{\pi}(s) = \mathbb{E}[G_t | s_t = s]$$

$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-t-1} R_T | s_t = s]$$

$$= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | s_t = s]$$

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

following policy π

Bellman equation for $q(s,a)$

$$q_{\pi}(s,a) = \mathbb{E}[G_t | s_t = s, A_t = a]$$

$$= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T | s_t = s, A_t = a]$$

$$= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | s_t = s, A_t = a]$$

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

following policy π

Solving a Markov decision process

The optimal value of a state is the expected return following the optimal policy

$$v_*(s) = \mathbb{E}_{\pi_*}[G_t | s_t = s]$$

$$q_*(s,a) = \mathbb{E}_{\pi_*}[G_t | s_t = s, A_t = a]$$

The optimal policy π_* is the one that chooses action that maximize $v_*(s)$, or $q_*(s,a)$

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

$$\pi_*(s) = \arg \max_a q_*(s,a)$$

To find the optimal policy π_* , we must know the optimal value

To find the optimal v_* or q_* value we must know the optimal policy

Bellman optimality equations

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

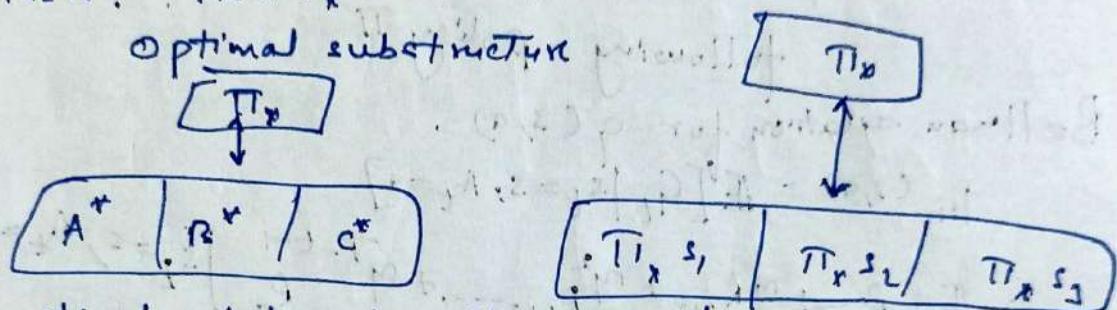
$$q_*(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} q_*(s',a')]$$

following the optimal policy π_*

Dynamic programming Method that finds the solution to a problem by breaking it down into smaller, easier problems

Our task: Find π_*

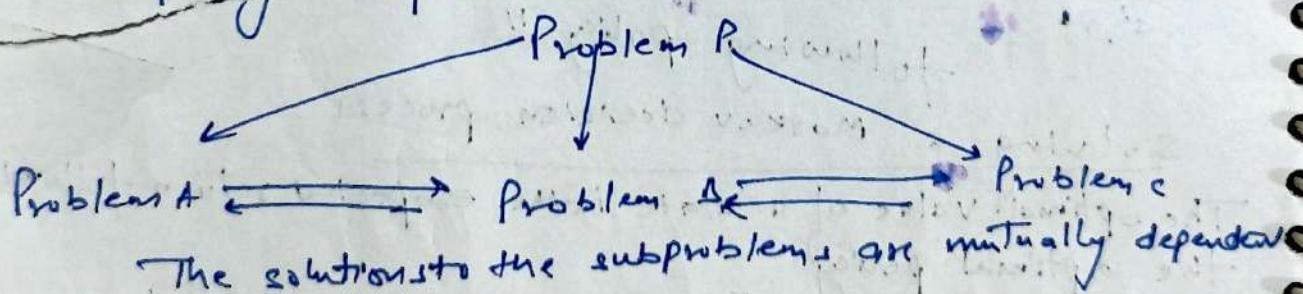
Optimal substructure



The optimal solution to all subproblems produce the optimal solution to the original problem.

$$\pi_* \iff \pi_*^*(s) \quad q_* \iff q_*^*(s,a) \quad v \iff V_*(s)$$

Overlapping subproblems



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

$$q_*^*(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} q_*(s',a')]$$

DP turns these eq. into update rules

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

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Model of the environment

$$V_s \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V_{\pi}(s')]$$

We need to know in advance how the state changes and what rewards we get from performing each action in each state
 DP: solves problems using expected values, not trial and error

Value iteration: We find the optimal policy π_* :

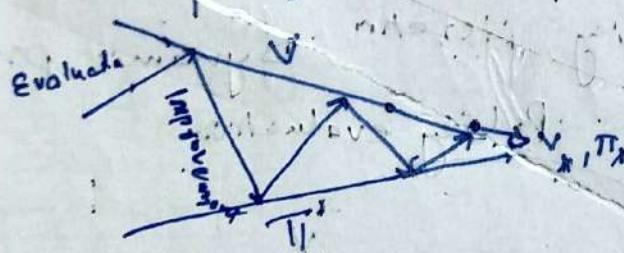
$$\pi_*(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V_*(s')]$$

... but for that you have to find V_* :

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Moves towards the value under the optimal policy

Policy iteration: A process that alternately improves the estimated value policy



Policy evaluation

Iterative policy evaluation

Iterative approximates the value of a given policy V_{π}

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

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

Each iteration gets closer to V_{π}

$$v_0 \rightarrow v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_{\pi}$$

Policy improvement: Does the policy improve if we change the first action?

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

If $q_{\pi}(s, \pi'(s)) > V_{\pi}(s)$ then

$$V_{\pi'}(s) \geq V_{\pi}(s)$$

π_1 and π_1' differ only in the action a they take in state s .

New policy! $\pi_1'(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi_1}(s')]$

Policy - stable = true

for $s \in S$ do

old action $\leftarrow \pi_1(s)$

$$\pi_1(s) \leftarrow \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi_1}(s')]$$

if old action $\neq \pi_1(s)$ then

policy - stable \leftarrow false

end if

end for

Policy iteration

Generalized policy iteration

Dynamic programming

Policy evaluation



Policy iteration results in the following iterative process:

then try to replicate the results of dynamic programming!

More efficiently

without a model

Disadvantage of dynamic programming

High computation cost

for each sweep we update all the states

complexity grows very rapidly with the no of states

Monte Carlo method family of methods that learn optimal $v_\pi(s)$ or $q_\pi(s, a)$ value based on samples

The agent will use a policy π to tackle the task for an entire episode

$$s_0, R_0, r_1, s_1, A_1, \dots, s_T, R_{T-1}, r_T$$

They approximate the values by interacting with the environment to generate sample returns and averaging them

$$v_\pi(s) = E_\pi [G_t | s_t = s] \quad G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$$

$$v_\pi(s) = \frac{1}{N} \sum_{k=1}^N G_{sk}$$

$$q_\pi(s, a) = E_\pi [G_t | s_t = s, A_{t-1} = a] \quad G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$$

$$q_\pi(s, a) = \frac{1}{N} \sum_{k=1}^N G_{sk, a}$$

Law of large numbers

In the limit, this succession of return samples

$G_{s1}, G_{s2}, \dots, G_{sn}$ converges to

$$P\left(\lim_{n \rightarrow \infty} \bar{G}_s = v_\pi(s)\right) = 1$$

Advantages of Monte Carlo Method

Dynamic Programming

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

$v_\pi(s)$ = max $\sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$ state value does not depend on the rest

- The estimate of a state value is independent of the total no. of states
- we can focus the estimations on the states that help us solve the task

$s_0, A_0, R_1, \dots, t_1, R_2, \dots, R_T$
 from the generated trajectory, we will calculate
 the returns for each moment of time t

$$q_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

$$q_{t+1} = R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots + \gamma^{T-t-1} R_T$$

our old strategy Use these returns to update the
 value function and through it, the policy π
 To improve the policy, we can use $v\pi$ argument

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

$v(s)$ requires knowing the effects of taking each
 action beforehand

~~To~~ The environment dynamics are implicit in v_π

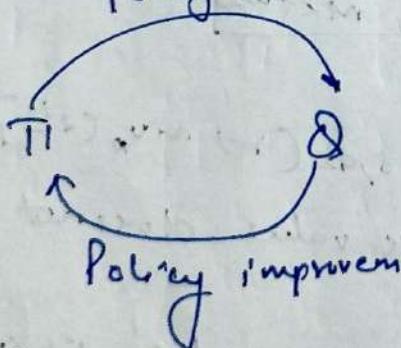
$$v_\pi(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$$

Instead of:

$v(s)$

we keep a table with:

Policy evaluation $Q(s, a)$



$$\pi_0 \rightarrow Q_{\pi_0} \rightarrow \pi_1 \rightarrow Q_{\pi_1} \rightarrow \dots$$

$$\dots \rightarrow Q_{\pi_K} \rightarrow \pi_K$$

The importance of exploration

$Q(s, a)$ is a estimate

$$\pi'(s) = \arg \max_a Q(s, a)$$

The estimate will improve as we
 obtain new samples, but the might
 not be perfect

If action a is optimal but has a bad estimate $Q(s, a)$
we'll never pick it

The only way to avoid this is to explore all actions
every once in a while and update their estimate
 $Q(s, a)$

How can we update maintain the exploration?

To option
exploring state starts
 $s \sim S, A_0 \sim \pi(s, \cdot)$ stochastic policies
 $\pi(a|s) > 0, \forall a \in A(s)$

stochastic policies

On-Policy learning

Generates samples using
the same policy π that
we're going to optimize

off-policy learning

Generates samples with
an exploratory policy π'
different from the one we're
going to optimize

On-Policy Monte Carlo

ϵ -greedy policy
with probability ϵ select a
random action

with probability $1 - \epsilon$ selected
the action with highest
 $Q(s, a)$

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \epsilon \gamma & a = a^* \\ \epsilon \gamma & a \neq a^* \end{cases} \quad \epsilon = \frac{\epsilon}{|A|}$$

$$|A| = 4, \epsilon = 0.2$$

$$\pi(a|s) = \begin{cases} 1 - 0.2 + 0.05 = 0.85 & a = a^* \\ 0.05 & a \neq a^* \end{cases}$$

$$\epsilon = \frac{0.2}{4} = 0.05$$

off policy Monte Carlo

off-policy strategy

Exploratory policy
b(a|s)

Target policy
 $\pi(a|s)$

off-policy strategy
b(a|s)

Generates the episode we are going to update $Q(s, a)$ with $s_0, A_0, R_1, s_1, A_1, \dots, R_T$
 $\pi(a|s)$

Policy to be optimized through $Q(s, a)$ values:
 $\pi(s) \leftarrow \arg \max_a Q(s, a)$

The policy b has to be able to explore all the actions that π it can take!

If $\pi(a|s) > 0$, then $b(a|s) > 0$.
The average return will not approximate the value under π but under b!

$$E_b [q_t | s_t = s, A_t = a] = q_b(s, a)$$

Importance sampling:- statistical technique for estimating the expected values of a distribution by working with samples from another distribution

$$w_t = \frac{\pi(A_t | s_t)}{b(A_t | s_t)}$$

By correcting the returns using importance sampling we will approximate the value under π !

$$E[\pi] E[w_t q_t | s_t = s] = v_\pi(s)$$

Update Rule

for each $Q(s, a)$, we'll keep a list of observed returns $[q_1, q_2, q_3, \dots, q_N]$

Each time we need to update $Q(s, a)$ we'll recompute the average

$$Q(s, a) \leftarrow \frac{1}{N} \sum_{k=1}^N G_k$$

update rule for $Q(s, a)$

$$Q(s, a) \leftarrow Q(s, a) + \frac{w_k}{C(s, a)} [G_k - Q(s, a)]$$

$$\text{where } C(s, a) = \sum_{k=1}^N w_k$$

Temporal difference methods

- Family of methods that learn the optimal $V_\pi(s)$ or $Q_\pi(s, a)$ values based on experience
- combination of Monte Carlo methods and dynamic programming

The agent learns from example

$$s_0, A_0, R_1, s_1, A_1, \dots, R_T$$

They do not need a model of the environment

$$T_1'(s) = \arg \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V(s')]$$

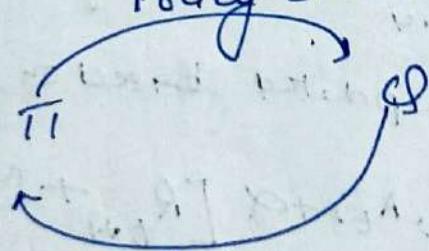
$$T_1'(s) = \arg \max_a Q(s, a)$$

$$\text{They use bootstrapping!} \quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

update rule for the STOCA algorithm

Generalize Policy Iteration

Policy Evaluation



Monte Carlo methods wait until the return G_t is available before updating $Q(s, a)$:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

we update Q and π every time the agent takes an action.

Temporal difference methods

we keep a table with q -value estimates for each s_t, A_t pair:

$$Q(s_t, A_t)$$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | s_t = s, A_t = a]$$

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

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

$$q_\pi(s, a) = \sum_{r, s'} p(r, s' | s, a) \left[r + \gamma \sum_{a'} \frac{\pi(a' | s')}{\pi(a)} q_\pi(s', a') \right]$$

Compare the two estimates:

$$R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(s_t, A_t)$$

Temporal difference errors.

Estimates are updated based on temporal difference error:

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(s_t, A_t)]$$

update rule for the SARSA method

Similar to constant α Monte carlo!

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_t + \gamma Q(s_{t+1}, a_{t+1})]$$

Estimating:

$$\hat{Q}_t = R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

Allows us to update $Q(s_t, a_t)$ at time $t+1$.

$$Q(s_t, a_t) \leftarrow (1-\alpha)(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})]$$

The update moves $Q(s_t, a_t)$ α percent in the direction of $R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$

Monte carlo vs temporal difference method

Monte carlo learning

SARSA on-policy temporal difference learning

ϵ -greedy policy

Probability ϵ

$a \sim A(s)$ random

Probability $1-\epsilon$

$a = \underset{a}{\operatorname{argmax}} Q(s, a)$

The name sarsa comes from the five values involved in the update rule:

$$s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})]$$

The ϵ -greedy policy π picks the next action, used to update $Q(s_t, a_t)$

Q learning off-policy temporal difference learning exploratory policy

Target Policy

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

$b(s)$

update rule:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, \pi(s_{t+1})) - Q(s_t, a_t)]$$

update $Q(s_t, a_t)$ biased on best available action.

$$A_{t+1} = \pi(s_{t+1}) = \arg\max_a Q(s_{t+1}, a)$$

Advantages Vs Monte Carlo

They allow us to update $Q(s, a)$ while experience is being collected.

In practice they converge faster.

Advantages Vs Dynamic programming

More efficient, focus the effort on the states that lead to goals.

Don't require a model of the environment.

n step temporal difference method

family of algorithms that learn based on experience using the n-step bootstrapping technique

update rule of the SARSA algorithm:

$$Q(s_t, a_t) \rightarrow Q(s_t, a_t) + \alpha \frac{R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)}{Q(s_t, a_t)}$$

Bootstrapping

$$R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

update target

1 step bootstrapping

$$R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

After performing an action, we replace the rewards by the $Q(s_{t+1}, a_{t+1})$ values

2-step bootstrapping:

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(s_{t+2}, A_{t+2})$$

2-step bootstrapping

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+2} + \gamma^3 Q(s_{t+2}, A_{t+2})$$

n-step bootstrapping

$$R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(s_{t+n}, A_{t+n})$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(s_{t+n}, A_{t+n})$$

n-step return estimate

$$Q(s_t, A_t) \rightarrow Q(s_t, A_t) + \alpha [G_{t:t+n} - Q(s_t, A_t)]$$

Since we need to observe n rewards, we need to wait until time $t+n$ to update $Q(s_t, A_t)$.

SARSA

$$G_{t:t+1} = R_{t+1} + \gamma Q(s_{t+1}, A_{t+1})$$

$$Q(s_t, A_t) \rightarrow Q(s_t, A_t) + \alpha [G_{t:t+1} - Q(s_t, A_t)]$$

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(s_{t+2}, A_{t+2})$$

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+2} + \gamma^3 Q(s_{t+2}, A_{t+2})$$

If $n \geq T$:

$$G_{t:t+n} = G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$$

$$Q(s_t, A_t) \rightarrow Q(s_t, A_t) + \alpha [G_{t:t+n} - Q(s_t, A_t)]$$

Monte carlo

$$Q_{t:t+n} = Q_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$$

$$Q(s_t, a_t) \rightarrow Q(s_t, a_t) + \alpha [G - Q(s_t, a_t)]$$

How does n affect the learning process?

$$Q_{t:t+1} = R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

$$Q_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$$



$$Q_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n Q(s_{t+n}, a_{t+n})$$

How does the choice of n affect the learning process?

Because $Q(s_{t+1}, a_{t+1})$ is an estimate of future rewards. The estimate improves throughout the learning process.

$$R_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

$Q(s_{t+1}, a_{t+1})$ introduces bias in the estimate.

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(s_{t+2}, a_{t+2})$$

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 Q(s_{t+3}, a_{t+3})$$

$$R_{t+1} + \gamma R_{t+2} + \gamma^{n-1} R_{t+n} + \gamma^n Q(s_{t+n}, a_{t+n})$$

The higher n , the more heavily discounted the estimate $Q(s_{t+n}, a_{t+n})$ by γ^n .

The higher n , the lower the bias.

Variance! Each reward is a random variable that depends on the state s and action a proceeding

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

- Even if they have the same expected value G_t and
- es will be very different from each other
- the higher the n , the higher the variance
- the higher the n , the lower the bias but the higher the variance

n -step SARSA

Combining SARSA with n -step bootstrapping

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} Q(s_{t+n}, a_{t+n})$$

update rule:

$$Q(s_t, a_t) \rightarrow Q(s_t, a_t) + \alpha [G_{t:t+n} - Q(s_t, a_t)]$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n Q(s_{t+n}, a_{t+n})$$

On-policy learning strategy with ϵ -greedy policy

Probability ϵ

Probability $1-\epsilon$

$a \sim \pi(s)$ randomly

$$a = \arg \max_a Q(s, a)$$

n-step SARSA in action

Continuous state spaces $S = (0, \infty)$
Solutions

Transforming the states
(section 7)

Using other algorithms
(sections 8-11)

Transforming the states

convert the continuous state space into a discrete state space

$$R \rightarrow \{s_0, s_1, s_2, \dots, s_N\}$$

state aggregation

Tile coding

state aggregation

- we went from having infinite states to having only 10

state aggregation

2D state aggregation, 1D state aggregation

Aggregate the state and generate a grid to having a single state

Tile coding: Average of the estimate

Creation of code tiling

Step 1: Divide each state dimension in segments.

Step 2: Expand or shrink the value range

Step 3: Discretize the segments.

Step 4: Repeat steps 1 to 3 n times.

function approximators

continuous state space

Disadvantages of state aggregation → Limited precision.
and tile coding → complexity $O(n^k)$

$n = \text{segments}$, $k = \text{dimension}$

We need precise alternative of limited complexity function approximators. (state value function)

How do we observe the value function?

The agent learns

based on experience

The function $V_*(s)$ and

$q_*(s, a)$ are not known

in advance

Policy evaluation

$T(s) \rightarrow V(s)$

Policy improvement

$$f_1(s/w) \rightarrow f_2(s/w) \rightarrow \dots \rightarrow f_n(s/w) \approx V(s)$$

Two examples

• state

$s = [s_1, s_2, \dots, s_n]$

Parameter vector

$$w = [w_1, w_2, w_3, \dots, w_n]$$

Linear approximation

$$f(s/w) \approx \hat{V}(s/w) = w \cdot s^T$$

$$f(s) = \hat{V}(s)$$

$$= w_1 s_1 + w_2 s_2 + \dots + w_n s_n$$

Inputs

$$\phi(s) = [s_1, s_1^2, \dots, s_1^K, \dots, s_n, \dots, s_n^K]$$

Parameter vector:

$$w = [w_0, w_1, w_2, \dots, w_n]$$

Polynomial approximator

$$f(s|w) \rightarrow \hat{V}(s|w)$$

$$= w \cdot \phi(s)^T$$

$$f(s) = \hat{V}(s) = w_0 + w_1 \cdot s_1^2 + w_2 \cdot s_1^4 + \dots$$

Advantages of using function approximators:

Efficiency $w = [w_0, w_1, w_2, \dots, w_n]$

They only require memory to save their parameters

Flexibility, Precision

Artificial neural network

They serve multiple purpose including function approximation.

Deep Sarsa

SARSA + Neural network

$$s = [s_1, s_2]$$

$$\hat{q} = [\hat{q}(s, a_1), \hat{q}(s, a_2)]$$

Neural network optimization

Mean sq. error

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N [y_i - \hat{y}_i]^2$$

We want to minimize the square of the errors of the neural network estimates

$$L(\theta) = \frac{1}{(K)} \sum_{i=1}^{(K)} \left[R_i + \gamma \hat{q}(s'_i, A'_i | \theta_{\text{target}}) - \hat{q}(s_i, A_i | \theta) \right]^2$$

Estimated

Target

$$y_i = R_i + \gamma \hat{q}(s'_i, A'_i | \theta_{\text{target}})$$

value towards which we want to push the estimate

$$\hat{y} = \hat{q}(s_i, A_i | \theta)$$

estimate of the value of a state-action pair

we calculate the gradient vector of the cost function with respect to the θ parameters

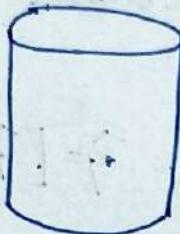
$$\nabla L(\theta) = \left[\frac{\partial L}{\partial \theta_0}, \frac{\partial L}{\partial \theta_1}, \dots, \frac{\partial L}{\partial \theta_n} \right]$$

with the gradient vector, we will make a SGD step

$$\theta \leftarrow \theta - \alpha \nabla L(\theta)$$

Replay memory

$s_0, A_0, R_0, s_1 \rightarrow$
 $s_1, A_1, R_1, s_2 \rightarrow$
 $s_2, A_2, R_2, s_3 \rightarrow$
...



Memory that stores the state transitions that the agent experiences

The memory has a limited size and when it fills up, it replaces old transitions with new ones

$$k = (s, t, R, s') \sim B$$

To update the NN, we randomly choose a batch of transitions from the memory

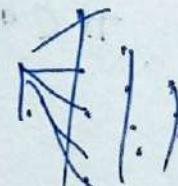
The batch of the transitions obtained from the memory is used to calculate the cost function and update the θ parameters

Target network

Boot strapping

function approximation

$$y_i = r_i + \gamma \hat{q}(s_i, a_i | \theta_{\text{target}})$$



when a value is changed, nearby values will also be affected

By modifying a $\hat{q}(s, a | \theta)$ estimate we also modify its $\hat{q}(s', a' | \theta_{\text{target}})$ target. we make a copy of the neural network to calculate the target

$$\theta_{\text{target}} \leftarrow \theta$$

This η does not change with SGD. Its θ parameters remain the same.

Each k episode, we copy the parameters from the agent network to the target network

$$\theta_{\text{tag}} \leftarrow \theta$$

Deep Q-learning

Q-Learning

Neural Network

$$S = [s_1, s_2]$$

$$\hat{q} = [\hat{q}(s, a_1), \hat{q}(s, a_2), \hat{q}(s, a_3)]$$

Experience replay

Policy gradient method
 Value based method
 function approximators
 Value table
 $a = \arg \max_a Q(s, a)$ $a = \arg \max_a \hat{q}(s, a | \theta)$

$$a_t = \arg \max_a Q(s, a)$$

$$q = \arg \max_q \hat{q}(s, a | \theta)$$

The policy is defined based on these values

Policy gradient method. :- They use a function approximator to estimate the probability of taking each action.

$$\pi_1(s|\varepsilon, \theta) \in [0,1]$$

Input

$$s = [s_1, s_2]$$

out put

$$\pi_1(s|\theta) = \{p(c_{q_1}), p(c_{q_2}), p(c_{q_3})\}$$

The man is the police

Advantages Value based method cannot represent the stochastic policy in a simple way.

Greedy policy

$$\pi_t(a'|s) = \begin{cases} 1 & \text{if } a' = \arg \max_a \hat{q}(s, a | \theta) \\ 0 & \text{else.} \end{cases}$$

Epsilon-greedy policy

$$\pi_t(a'|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A|} & \text{if } a' = \arg \max_a \hat{q}(s, a | \theta) \\ \frac{\epsilon}{|A|} & \text{else} \end{cases}$$

Si $\pi^*(s) = [0.7, 0.3]$, How do we represent it?

The policy changes more smoothly during learning.

Value based methods:
when the maximum q-value changes
they choose a new action 100% of
the time $a = \arg \max_a \hat{q}(s, a)$

Policy gradient method

The probability of choosing an
action changes in small increments
 $a \sim \pi(s | \theta)$

The Bellman Equation

concept

s - state

a - Action

R - Reward

γ - Discount

$$V(s) = \max_a (R(s,a) + \gamma V(s'))$$

Markov decision Process (MDP)

Deterministic search

↓
100%

Non-deterministic
search

↓
10% / 10%
80% ↗

$$0.8 * V(s_1') + 0.1 * V(s_2') + 0.1 * V(s_3')$$

$$V(s) = \max_a (R(s,a) + \gamma \overbrace{V(s')}^{\text{next}})$$

$$V(s) = \max_a (R(s,a) + \gamma \sum_{s'} P(s,a,s') V(s'))$$

Policy and vs Play

Adding a Living Penalty

Q Learning Intuition

$$V(s) = \max_a (R(s,a) + \gamma \sum_{s'} P(s,a,s') V(s'))$$

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} [P(s,a,s') V(s')] \downarrow \\ \max_{a'} Q(s',a')$$

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} [P(s,a,s') \max_{a'} Q(s',a')]$$

Temporal difference:

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} [P(s,a,s') \max_{a'} Q(s',a')]$$

$$Q(s,a) = R(s,a) + \gamma \max_{a'} Q(s',a')$$

Before

$$Q(s,a)$$

After

$$R(s,a) + \gamma \max_{a'} Q(s',a')$$

$$TDQ(s,a) = R(s,a) + \gamma \max_{a'} Q(s',a') - Q(s,a)$$

$$Q(s,a) = Q(s,a) + \alpha TDQ(s,a)$$

$$\Phi_t(s, a) = \Phi_{t-1}(s, a) + \alpha T D_t(a, s)$$

$$\Phi_t(s, a) = \Phi_{t-1}(s, a) + \alpha (R(s, a) + \gamma \max_{a'} \Phi(s', a') - \Phi_{t-1}(s, a))$$

RL I Multiarm Bandits

very close to real not exactly same to real returns etc

ϵ -Greedy approach

Reward, Action, Value of action (Φ)

$$\Phi_{n+1} = \Phi_n + \frac{1}{n} (R_n - \Phi_n)$$

$$\Phi_{n+1} = \Phi_n + \alpha (R_n - \Phi_n)$$

$$= \alpha R_n + (1-\alpha) \Phi_n$$

$$= \alpha R_n + (1-\alpha) [\Phi_{n-1} + \alpha (R_{n-1} - \Phi_{n-1})]$$

$$= \alpha R_n + (1-\alpha) [\alpha R_{n-1} + (1-\alpha) \Phi_{n-1}]$$

$$= \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \Phi_{n-1}$$

$$= (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha(1-\alpha)^{n-i-1} R_i$$

Noise strategy

1. Explore only

$$\begin{array}{ccc}
 100 & 100 & 100 \\
 \times 10 & \times p & \times 5 \\
 \hline
 1000 & 100p & 500 \\
 \hline
 & \underbrace{\qquad\qquad\qquad}_{\text{Sum } 2200} &
 \end{array}$$

$$[P = 700]$$

2. Exploit only

$$R_1 = 7 \quad R_2 \rightarrow 8 \quad R_3 \rightarrow 5$$

$$20 + 297(8) = 2396$$

$$[P = 100]$$

3. ϵ -Greedy

$$\epsilon = 10\%$$

↳ 30 Explore

↳ 270 Exploit

Result 2907 Avg.

$$[P \approx 100]$$

Zero regret strategy

$$\frac{P}{T} \rightarrow 0$$

$$T \rightarrow \infty$$

$\hat{Q}_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken}}{\text{prior } \hat{n}_t}$

number of times a taken
prior to t

$$= \frac{\sum_{i=1}^{t-1} R_i \cdot 1_{A_i=a}}{\sum_{i=1}^{t-1} 1_{A_i=a}}$$

Greedy action selection rule

$$A_t \doteq \underset{a}{\operatorname{argmax}} \hat{Q}_t(a)$$

New Estimate \leftarrow old estimate +
step size [Target - Old estimate]

$$\begin{aligned}\hat{Q}_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\ &= \frac{1}{n} \left(R_n + \sum_{i=1}^{n-1} R_i \right)\end{aligned}$$

$$= \frac{1}{n} \left(R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right)$$

$$= \frac{1}{n} \left(R_n + (n-1) \hat{Q}_n \right)$$

$$= \frac{1}{n} (R_n + nQ_n - Q_n)$$

$$= Q_n + \frac{1}{n} [R_n - Q_n]$$

Exponential Recency - Weighted Average

$$Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$

$$= \alpha R_n + (1-\alpha) Q_n$$

$$= \alpha R_n + (1-\alpha) [\alpha R_{n-1} + (1-\alpha) Q_{n-1}]$$

$$= \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$= \alpha R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2} +$$

$$\dots + (1-\alpha)^{n-1} \alpha R_1 + (1-\alpha)^n Q_0$$

$$= (1-\alpha)^n Q_0 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

UCB

$$A_t = \arg \max_a \left[Q_t(a) + C \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Gradient Bandit Algorithm

$$\Pr \{ A_t = a \} = \frac{e^{H_t(a)}}{\sum_{b=1}^K e^{H_t(b)}} = \pi_t(a)$$

Gradient Bandit. Update with
stochastic Gradient Ascent

$$H_{t+1}(A_t) = H_t(A_t) + \alpha(R_t - \bar{R}_t) \\ (1 - \pi_t(A_t))$$

$$H_{t+1}(a) = H_t(a) - \alpha(R_t - \bar{R}_t) \pi_t(a)$$

Upper confidence Bound
(UCB)

hope around estimate value

$$Q_A + \sqrt{\frac{2 \ln N}{n_A}}$$

Qvalue :- simply average reward from machine

$$Q_B + \sqrt{\frac{2 \ln N}{n_B}}$$

n_A = no. of times machine A is played

$$Q_C + \sqrt{\frac{2 \ln N}{n_C}}$$

N = Total plays in all the machines

$$so \quad n = n_A + n_B + n_C$$

Thompson sampling

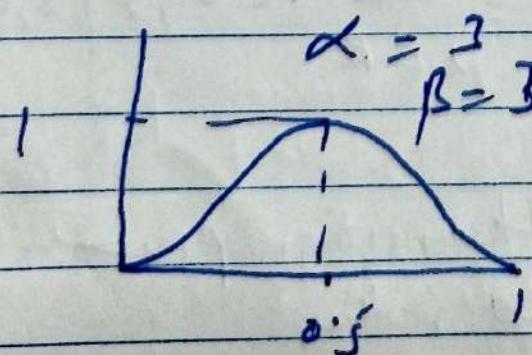
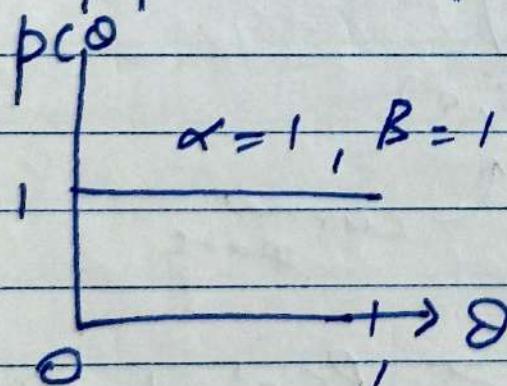
B-Distribution

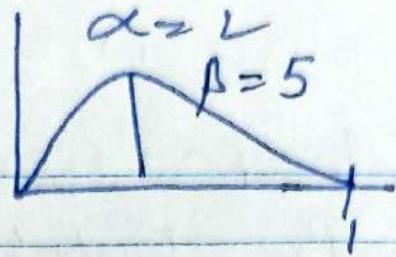
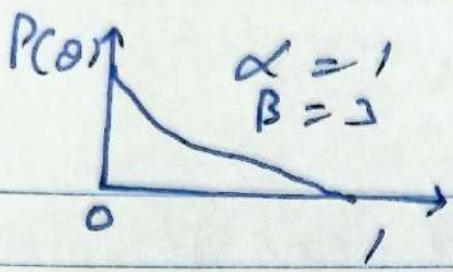
$$p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\theta \in [0, 1]$$

where, $\Gamma(\cdot)$ gamma function

$\alpha, \beta \rightarrow$ shape parameters > 0





Markov Decision Process

1. Agent
2. Environment
3. Reward
4. state

Policy iteration value iteration
Policy $P(s'|s)$

Transition dynamics

$$P(s', r | s, a)$$

(next state current state)

$$P(s' | s, a) = \sum_{r \in R} P(s', r | s, a)$$

state transition probability

Total discount return

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{n-1} R_{t+n}$$

Value function:

(with respect to a policy π)

$$V_\pi(s) = \mathbb{E}_\pi [q_t | s_t = s]$$

$$V_\pi(s) \geq \mathbb{E}_\pi [q_t | s_t = s]$$

$$= \mathbb{E}_\pi [r_{t+1} + \gamma q_{t+1} | s_t = s]$$

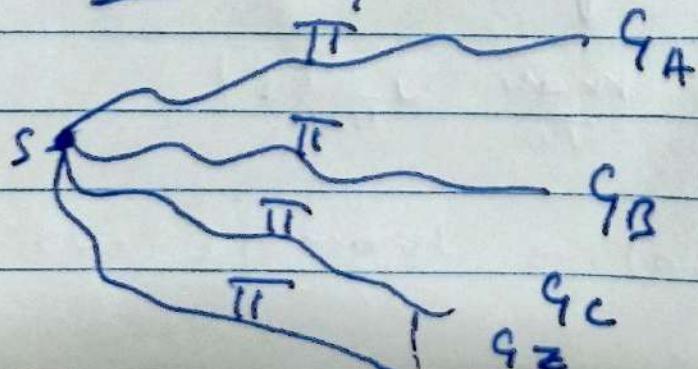
$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a)$$

current state $\underbrace{\left[r + \gamma \mathbb{E}_\pi [q_{t+1} | s_{t+1} = s] \right]}_{\text{policy dynamics}}$

$$\Rightarrow V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a)$$

$$\left[r + \gamma V_\pi(s') \right]$$

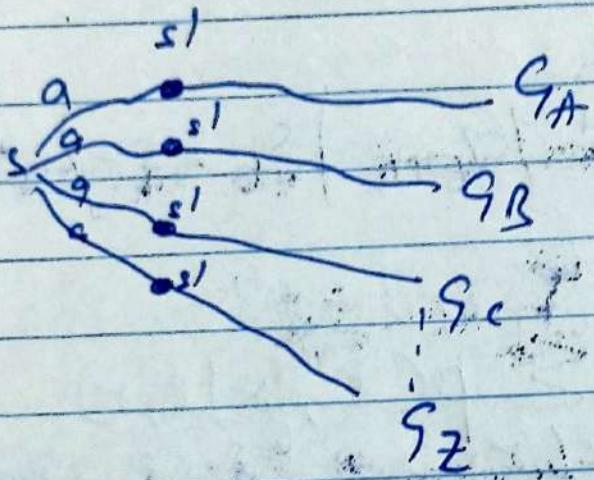
Bellman equation next states value



$$V_{\pi}(s) = \frac{1}{26} \sum (G_A + G_B + \dots + G_Z)$$

Value of an action:

$$q_{\pi}(s, a) = E_{\pi} [G_t | s_t = s, A_t = a]$$



$$\hat{q}_{\pi}(s, a) = \frac{1}{26} \sum (G_A + G_B + \dots + G_Z)$$

optimal state value and optimal action value

optimal \rightarrow Maximum possible value of a state / action.

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

UCB derivation

$$L_t = \sum_a n_t(a) \Delta_a$$

$$\Delta_{a^*} = 0$$

$$\forall a = a^* \quad n_t(a) \Delta_a \leq x_a \log t$$

$$\boxed{A_t = \arg\max_a Q_t(a) + U_t(a)}$$

$$m \leq t \quad n_m(a) \Delta_a \leq x_a \log m \leq x_a \log t$$

$$n \in \{m+1, \dots, t\}: n_n(a) \Delta_a > x_a \log n$$

$$\mathbb{E}[n_t(a)] = \mathbb{E}\left[\sum_{n=1}^t \mathbb{I}(A_n=a)\right]$$

$$= \mathbb{E}\left[n_m(a) + \sum_{n=m+1}^t \mathbb{I}(A_n=a)\right]$$

$$= \mathbb{E}\left[x_a \frac{\log t}{\Delta_a} + \sum_{n=m+1}^t \mathbb{I}(A_n=a)\right]$$

$$= x_a \frac{\log t}{\Delta_a} + \sum_{n=m+1}^t \mathbb{E}\left[\mathbb{I}(A_n=a) \middle| N_n(a) \geq x_a \log n\right]$$

$$P(A_n = a) \leq [N_n(a) \Delta_a]^{-x_a \log n}$$

$P(A_n = a) ?$ Given $N_n(a) \Delta_a > x_a \log n$

$$\leq P(\phi_t(a) + u_t(a) \geq \phi(a^*) + u_t(a^*))$$

Simplified the equation

$$= P(\phi + u \geq \phi^* + u^*) ? = 1$$

$$= P(\phi + u \geq \phi^* + u^* \mid \underbrace{\phi^* + u^* \leq y}_{P(\phi^* + u^* \leq y)})$$

$$+ P(\phi + u > \phi^* + u^* \mid \underbrace{\phi^* + u^* > y}_{P(\phi^* + u^* > y)})$$

$$\frac{1}{n}$$

$$y \stackrel{?}{=} q(a^*)$$

$$q(a^*)$$

$$P(\phi^* + u^* \leq q^*) \leq e^{-N_n(a^*) u_n(a^*)^2} = \frac{1}{n}$$

$$\sum_{n=1}^t \frac{1}{n} \leq \log t + 1$$

$$e^{-N u^2} = \frac{1}{n} \Rightarrow u = \sqrt{\frac{\log n}{2}}$$

$$P(C\varphi + u \geq \varphi^* + u^*) \leq \frac{1}{n} \quad | \quad (\varphi^* + u^* > q^*)$$

$$\leq P(C\vartheta + u \geq q^*)$$

$$= P(C\vartheta - u \leq -q^* + q - q) \\ = -\Delta_a$$

$$= P(C\vartheta - \Delta_a - u \leq q)$$

$$\Delta_a \geq 2u ?$$

u?

$$n \cdot \Delta_a > x_a \log n$$

$$\Delta_a > x_a \frac{\log n}{N}$$

$$x_a \leq t(a)^2$$

$$x_a = \frac{4}{\Delta_a} \Rightarrow \Delta_a^2 > 4t(a)^2$$

$$\Rightarrow \Delta_a > 2u$$

$$\leq PC - \varphi + 2u - 4 \leq -9$$

$$= PC - \varphi + u \leq -9 \quad | \quad u = \sqrt{\frac{\log n}{n}}$$

$$\leq e^{-\frac{1}{n}} = \frac{1}{n}$$

Transfer learning Deep learning
All weights in mind

- used in mostly in CNN vision
- " Alex net 8 layers
Imagenet 1M, 1000 class

→ 1000 image

→ Ordinary system

Less data

Down up
deep learning

use transfer learning

- enabling mechanism

- weights same

weights and architecture same