

2012 MARCH

APRIL 2012						
WEEK	M	T	W	T	F	S
13/18	30	1	2	3	4	5
14	6	7	8	9	10	11
15	12	13	14	15	16	17
16	18	19	20	21	22	23
17	25	26	27	28	29	

MODULE 1.

Monday

12

DAY 072-294 WK 11

MEETINGS / APPOINTMENTS

agent \Rightarrow decision maker

\hookrightarrow chooses between actions (k -actions)

based on the choice gets reward

action-values :- expected ~~value~~ reward of that action.

$$q_{V_\pi}(a) := \mathbb{E}[R_t | A_t = a] \quad \forall a \in \{1, 2, \dots, k\}$$

↓
 action value
 of action a expectation
 ↓
 Reward. Action = a

for each possible action.

$$q_{V_\pi}(a) = \sum_r p(r|a) r$$

↓
probability of observing that reward

possible reward

Maximize this:-

$$\underset{a}{\operatorname{argmax}} q_{V_\pi}(a)$$

\hookrightarrow find a .

IMPORTANT

13

Tuesday

WK 11 DAY 073-293

MARCH 2012

MARCH 2012						
WEEK	M	T	W	T	F	S
9				1	2	3
10	5	6	7	8	9	10
11	12	13	14	15	16	17
12	19	20	21	22	23	24
13	26	27	28	29	30	31

MEETINGS / APPOINTMENTS

2. Estimating Action Values

Sample-average method } → Estimate action values

9

Action Value : $a_t(a)$

10

$Q_t(a) := \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{no. of times } a \text{ taken prior to } t}$

11

12

$$Q_t(a) = \frac{\sum_{i=1}^{t-1} R_i}{t-1}$$

13

value of action at time t depends on rewards of actions taken before time t .

14

Greedy action selection → selecting the action that currently has the largest estimated value.

15

Non-greedy action → sacrifice immediate reward hoping to get more info about other actions

IMPORTANT

2012 MARCH

2012

APRIL	M	T	W	T	F	S
WK	30	1	2	3	4	5
13/16	6	7	8	9	10	11
14	12	13	14	15	16	17
15	18	19	20	21	22	23
16	25	26	27	28	29	24
17	23	24	25	26	27	28

Wednesday

14

DAY 074-292 WK 11

3. Exploration vs Exploitation

Exploration vs Exploitation.

non-greedy greedy

Incremental Implementation

↳ for not having to store millions of data.

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

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

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

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

new = old estimate + step size (Target - old estimate)

$$Q_{n+1} = Q_n + \alpha_n (R_n - Q_n)$$

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

$$\alpha \approx \frac{1}{n}$$

IMPORTANT

MARCH 2012

15

Thursday

WK 11 DAY 075-291

MEETINGS / APPOINTMENTS

WK	M	T	W	T	F	S	S
9					1	2	3
10	5	6	7	8	9	10	11
11	12	13	14	15	16	17	18
12	19	20	21	22	23	24	25
13	26	27	28	29	30	31	

→ Solution: constant
step size.

Non-stationary } Distribution of Rewards
Bandit Problem } changes with time.

9

If $\alpha_n = \text{Step size}$ → fixed to say 0.1
then recent rewards affect more than older rewards.

$$1 Q_{n+1} = Q_n + \alpha_n (R_n - Q_n)$$

$$2 Q_n = \alpha_n R_n + (1-\alpha) Q_{n-1}$$

$$3 = \alpha_n R_n + (1-\alpha) \alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2}$$

$$4 + \dots + (1-\alpha)^{n-1} \alpha R_1 + (1-\alpha)^n Q_1$$

$$5 Q_{n+1} = (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

initial action value weighted sum of the rewards at all the times

$(1-\alpha)^n Q_1$: contribution of Q_1 decreases exponentially with time.

$\sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$: the older rewards contribute exponentially less to the sum.

2012 MARCH

2012

APRIL	M	T	W	F	S	S
13/18					1	
14	30	2	3	4	5	6
15		7	8	9	10	11
16		12	13	14	15	16
17		18	19	20	21	22
	23	24	25	26	27	28
	29					

Friday

16

DAY 076-290 WK 11

MEETINGS / APPOINTMENTS

Exploration vs Exploitation trade off

↓
improve for
long-term
benefit.

↓
exploits knowledge
for short-term benefit

↳ eventually we'll have
more accurate estimates of action values.

To choose between explore & exploit - - -

Epsilon-Greedy Action Selection

↳ Greedy Exploit with a small probability ϵ

↳ Explore with a small probability

↳ Exploit otherwise, i.e., higher probability

$A_t \leftarrow \begin{cases} \text{argmax}_a Q_t(a) & \text{with prob. } 1-\epsilon \\ a \sim \text{Uniform}(\{a_1, a_2, \dots, a_k\}) & \text{with prob. } \epsilon \end{cases}$

$\text{argmax}_a Q_t(a) :=$ choose a with $\max_a Q_t(a)$ Greedy

$a \sim \text{Uniform}(\{a_1, a_2, \dots, a_k\}) :=$ choose randomly
bet. a_1, a_2, \dots, a_k
→ Explore

IMPORTANT

17

Saturday

WK 11 DAY 077-289

MARCH 2012

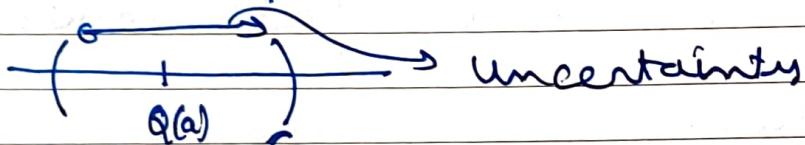
MARCH 2012						
WEEK	M	T	W	T	F	S
9					1	2
10	5	6	7	8	9	10
11	12	13	14	15	16	17
12	19	20	21	22	23	24
13	26	27	28	29	30	31

Optimistic Initial Values

→ Initially assume reward to be 1. or even more.

Upper - Confidence Bound (UCB)

↳ uses uncertainty in estimates to drive exploration.



Estimated value

upper bound \equiv highest possible reward of this action.

↳ choose the action with highest upper bound. (highest possible reward)

18

Sunday

WK 11 DAY 078-288

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

action that has
(highest estimated)
value. \equiv exploration
+ UCB
 \equiv exploitation. \equiv exploration
term

$c \rightarrow$ user defined: controls amt. of exploration.
 \equiv exploration. Eg: $c = 2$

IMPORTANT

2012 MARCH

APRIL 2012						
M	T	W	T	F	S	S
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29						

Monday

19

DAY 079-287 WK 12

MEETINGS / APPOINTMENTS

$t = \text{time steps}$; $N_t(a) = \text{times action } a \text{ taken}$

9

→ initially UCB explores more
UCB's exploration reduces over time.

10

→ ϵ -greedy always explores with ϵ % probability

12

Contextual Bandits :-

1

- Repeatedly :
1. Observe features x
 2. Choose action $a \in A$
 3. Observe reward r

2

Goal : Maximize reward.

3

4

5

6

APR

MAY

JUN

IMPORTANT

20

Tuesday

MARCH 2012

WK 12 DAY 080-286

MODULE: 02

MEETINGS / APPOINTMENTS

1. Markov Decision Process (MDP)

9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

It takes into account

different situations call for different actions.

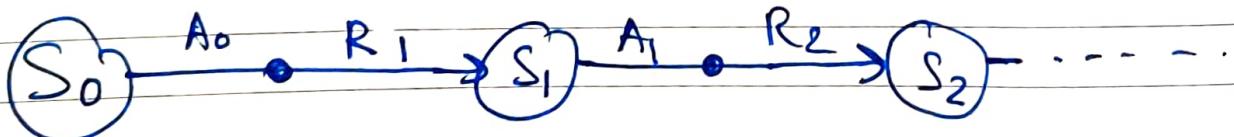
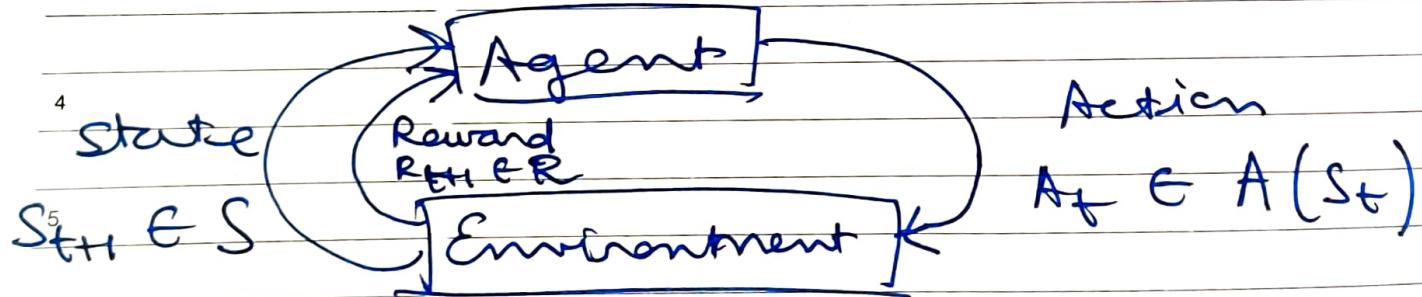
6 K-armed bandit does not

7 It takes into account

long term impact of our decisions

8 K-armed bandit : immediate rewards

2 State \equiv situation in which one action is taken.



Actions influence [immediate Rewards] & [future States]

↓
Future Rewards

IMPORTANT

2012 MARCH

APRIL 2012						
M	T	W	T	F	S	S
1						
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29

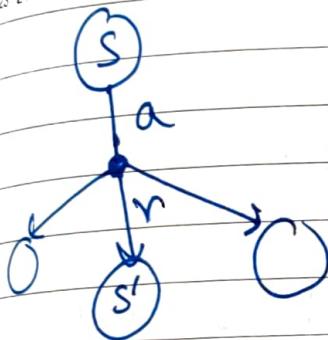
Wednesday

21

DAY 081-285 WK 12

MEETINGS / APPOINTMENTS

Dynamic
fnc. p
 $p(s', r | s, a)$



$$p: S \times R \times S \times A \rightarrow [0, 1] \quad ^9$$

$$\sum_{s' \in S} \sum_{r \in R} p(s, r | s, a) = 1, \forall s \in S \quad ^{10}$$

$$a \in A \quad (s)$$

(Future state & reward) depend on (Present state & action) Markov property¹²

Present state is sufficient. Older states are not that imp. ^2

Example : Task : The goal of the robot is to pick and place objects. ^3

State : latest readings of joint angles & velocities ^4

Action : amt. of voltage applied to each motor. ^5

Reward : +100 when an object is successfully placed. ^6

- 1 for each unit of energy consumed.

IMPORTANT

22

Thursday

WK 12 DAY 082-284

MARCH 2012

MARCH 2012					
W	K	M	T	W	F
9				1	2 3 4
10	5	6	7	8 9 10	11
11	12	13	14	15 16 17	18
12	19	20	21	22 23 24	25
13	26	27	28	29 30	31

MEETINGS / APPOINTMENTS

2. Goal of RL

Goal of an agent

$$G_t := R_{t+1} + R_{t+2} + R_{t+3} + \dots$$

9

9

10

return. total future reward \Rightarrow maximize

11

Dynamics of the MDP can be stochastic

12

$$\mathbb{E}[G_t] = \mathbb{E}[R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T]$$

1

2

expected.

interaction ends

[Maximize expected ~~return~~ return]

3

Episodic Tasks

→ Break into episodes

→ Each episode begins independently of how the previous one ended.

→ Every episode has a final state: terminal state.

After terminal state, each episode is reset.

2012 MARCH

APRIL 2012						
WEEK	M	T	W	T	F	S
13/18	30	1	2	3	4	5
14	6	7	8	9	10	11
15	12	13	14	15	16	17
16	18	19	20	21	22	23
17	24	25	26	27	28	29

Friday

23

DAY 083-283 WK 12

MEETINGS / APPOINTMENTS

Reward Hypothesis.

1. Give a man a fish & he'll eat for a day.

PROGRAMMING (noFAI)

9

2. Teach a man to fish & he'll eat for a lifetime.

SUPERVISED LEARNING (LfD)

11

3. Give a man a taste for fish and he'll figure out how to get fish, even if the details change.

OPTIMIZATION (RL)

12

⇒ Using RL to reposition solar energy

Common currency or unit of cost

↳ Moving the motors consumes energy
↳ Sun shining directly brings in energy

4

5

Goal Rewarding Action-Penalty.

↳ Goal : 1

↳ Not Goal : - 1

↳ Otherwise : 0

↳ Not Goal : 0

6

IMPORTANT

MARCH 2012

24

Saturday

WK 12 DAY 084-282

MARCH

2012

Wk	M	T	W	F	S
9				1	2 3 4
10	5	6	7	8	9 10 11
11	12	13	14	15	16 17 18
12	19	20	21	22	23 24 25
13	26	27	28	29	30 31

3. Continuing Tasks

Continuing Tasks vs Episodic Tasks

- Interaction goes on continuously.
- No terminal state.
- Interaction breaks naturally into episodes
- Each episode ends with a terminal state

$$G_t := R_{t+1} + R_{t+2} + \dots \stackrel{?}{=} \infty$$

↳ need to modify this
to make finite
↓

→ Episodes are independent

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

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

$0 \leq \gamma < 1$: Discount

25

Sunday

immediate rewards contribute more.

WK 12 DAY 085-281

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

$$G_t \leq \sum_{k=0}^{\infty} \gamma^k R_{\max} = R_{\max} \sum \gamma^k = \frac{R_{\max}}{1-\gamma}$$

$$G_t \leq \frac{R_{\max}}{1-\gamma} \Rightarrow \text{finite.}$$

IMPORTANT

2012 MARCH

APRIL 2012						
WEEK	M	T	W	T	F	S
13/18	30	1	2	3	4	5
14	6	7	8	9	10	11
15	12	13	14	15	16	17
16	18	19	20	21	22	23
17	24	25	26	27	28	29

26
Monday

DAY 086-280 WK 13

MEETINGS / APPOINTMENTS

 $\gamma = 0 \Rightarrow$ short-sighted agent $\gamma \rightarrow 1 \Rightarrow$ far-sighted agent

$$G_t = R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots)$$

$$\boxed{G_t = R_{t+1} + \gamma G_{t+1}}$$

Return can be defined recursively.

Specifying Policies

Policies : distribution over actions for each possible state.

Deterministic Policy

↳ Policy maps each state to a single action.

$$\Pi(s) = a$$

↳ agent can select same action in multiple states and some actions may be never selected

↳ Table representation.

MM
IMPORTANT

MARCH 2012

27

Tuesday

WK 13 DAY 087-279

MEETINGS / APPOINTMENTS

MARCH						
WEEK	M	T	W	T	F	S
9				1	2	3
10	5	6	7	8	9	10
11	12	13	14	15	16	17
12	19	20	21	22	23	24
13	26	27	28	29	30	31

Stochastic policy

Policy assigns probabilities to each action in each state.

$\pi(a|s)$: probability of selecting an action a in a state s .

$$\sum_{a \in A(s)} \pi(a|s) = 1 \quad \& \quad \pi(a|s) \geq 0$$

State has all the info required for decision making

↳ action depends only on current state & nothing else.

IMPORTANT

2012 MARCH 2012 value of a state, under a policy π , is the expected return from that state, if we act according to π . ↳ This is a discounted sum of future rewards.

28 Wednesday

APRIL	M	T	W	T	F	S	S
WK	1						
13/18	30	1	2	3	4	5	6
14		7	8	9	10	11	12
15		13	14	15	16	17	18
16		19	20	21	22	23	24
17		25	26	27	28	29	

Delayed Reward

MODULE : 03

MEETINGS / APPOINTMENTS

1. Value Functions

$$V(s) := \mathbb{E} [G_t | S_t = s] ; G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

10

State value func. Is the future reward an agent can expect to receive, starting from a particular state.

11

↳ Expected return from a given state

$$V_\pi(s) := \mathbb{E}_\pi [G_t | S_t = s] \xrightarrow[\text{w.r.t. policy } \pi]{\substack{\text{state} \\ \text{value func.}}}$$

1

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

3

Action value func. of a state is the expected return of the agent selection action a and then follows policy π .

4

Value functions predict rewards into the future.

5

→ State-value func. represent the expected return from a given state under a specific policy.

6

→ Action-value func. represent the expected return from a given state after taking a specific action, later following a specific policy.

IMPORTANT

APR

MAY

JUN

MARCH 2012

29

Thursday

WK 13 DAY 089-277

MARCH

WK	M	T	W	T	F	S	S
9				1	2	3	4
10	5	6	7	8	9	10	11
11	12	13	14	15	16	17	18
12	19	20	21	22	23	24	25
13	26	27	28	29	30	31	

MEETINGS / APPOINTMENTS

2 Bellman Equation.

Relate the value of the current state to the value of the future state, without waiting to observe all the future rewards.

10

Bellman equation for state-value function

11

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s] \quad \left. \begin{array}{l} \text{state value} \\ \text{function} \end{array} \right\}$$

12

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

$$2 \quad V_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$3 \quad V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

4 Defines a relationship between the value

5

of a state: $V_{\pi}(s)$ & the value of its6 possible successor states: $V_{\pi}(s')$

Bellman equation for action-value function

$$7 \quad Q_{\pi}(s, a) := \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

8 In this case the eqn does not begin with the policy selecting an action: $\sum_a \pi(a|s)$

IMPORTANT

2012 MARCH

	APRIL 2012					
WK	M	T	W	T	F	S
13/18	30					1
14	2	3	4	5	6	7
15	9	10	11	12	13	14
16	16	17	18	19	20	21
17	23	24	25	26	27	28

Friday

30

DAY 090-276 WK 13

MEETINGS / APPOINTMENTS

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

there is also a wt.

of probability under π of selecting a' in the state s'

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

→ Use the Bellman eqⁿ to compute value func by writing a system of linear equations.

→ We can only solve small MDPs directly, but Bellman eqns will factor into the solutions we see later for large MDPs.

p := Dynamic Func. ⇒ May be stochastic.

s' := all immediate next step states possible.

IMPORTANT

MARCH 2012

31

Saturday

WK 13 DAY 091-275

MARCH

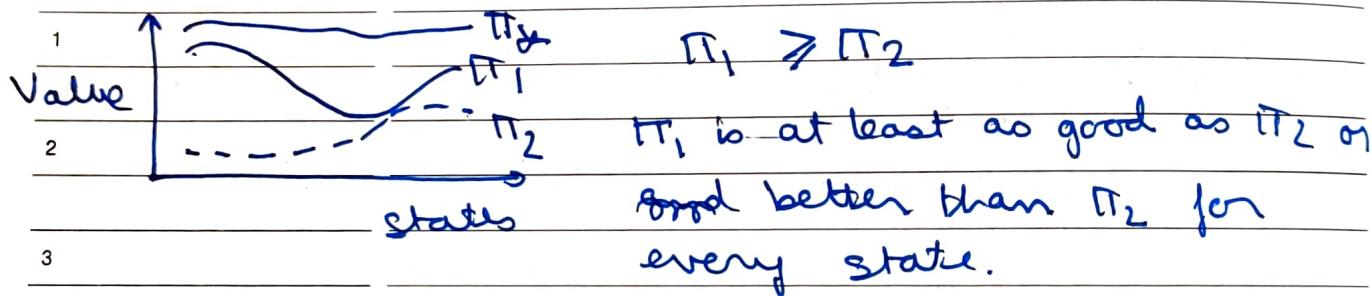
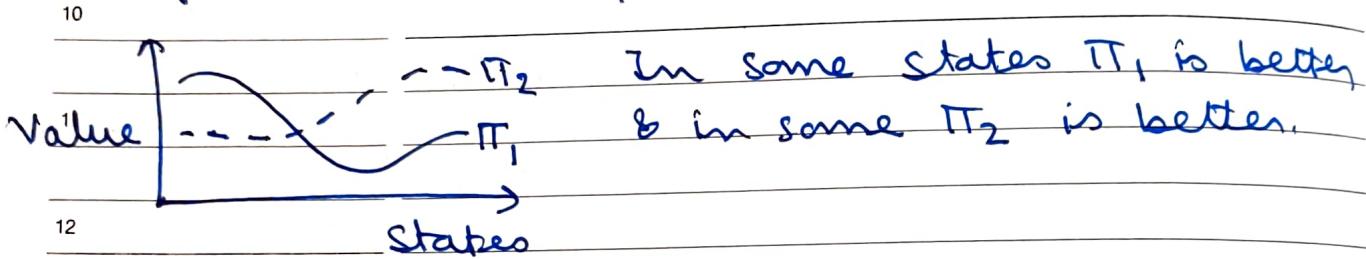
WK	M	T	W	T	F	S
9				1	2	3
10	5	6	7	8	9	4
11	12	13	14	15	16	11
12	19	20	21	22	23	18
13	26	27	28	29	30	25

MEETINGS / APPOINTMENTS

B - Optimal Policies

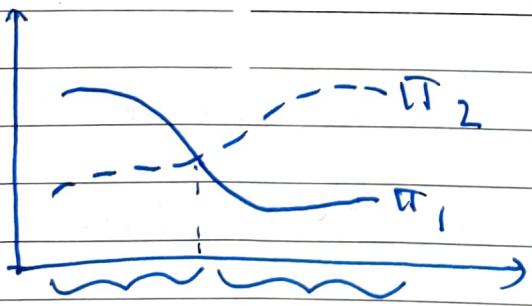
The policy specifies how an agent behaves.

Given this way of behaving, we then aim to find the value function.



An optimal policy π^* is as good as or better than all the other policies

$\Rightarrow \pi^*$ will have highest value in all states.



$$\text{TOP } \pi_3(s) = \pi_2(s)$$

$$\pi_3(s) = \pi_1(s)$$

$\pi^* \equiv \pi_3 \rightarrow$ a third policy which is the policy with the highest value in the current state.

IMPORTANT

01

Sunday

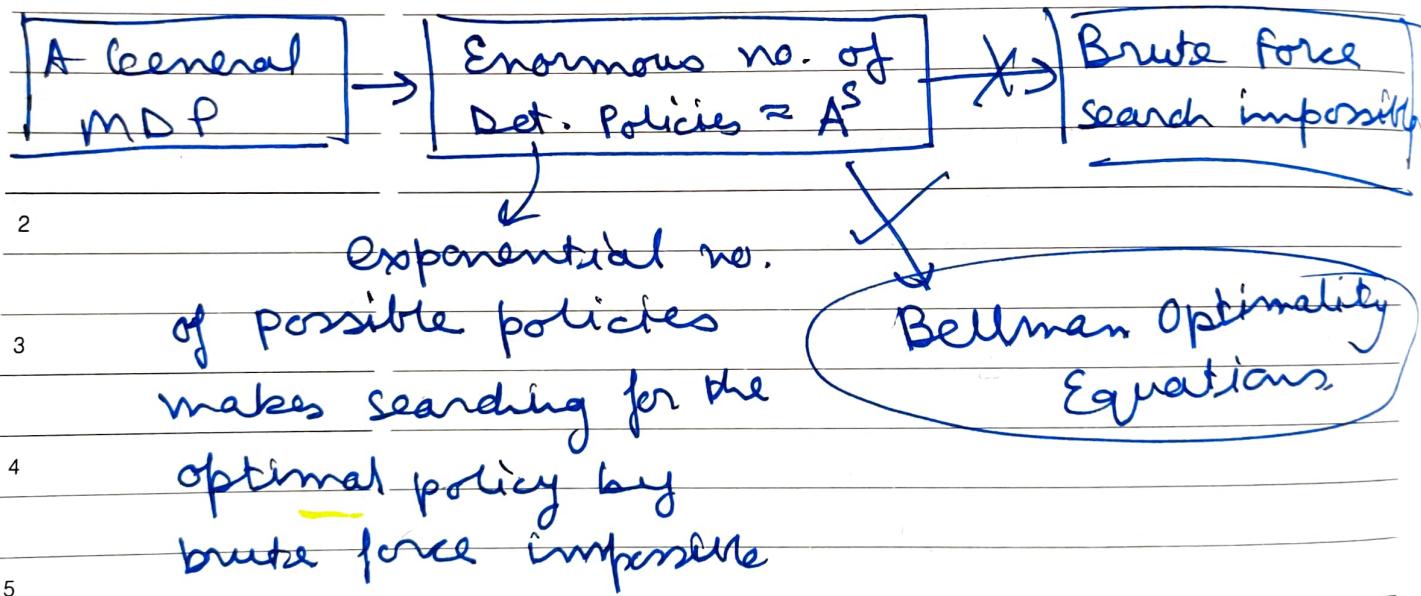
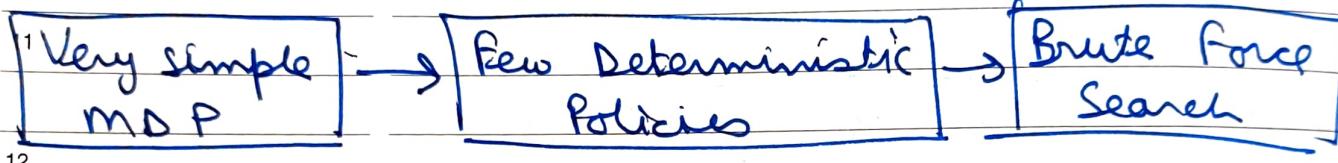
WK 13 DAY 092-274

MEETINGS / APPOINTMENTS

APRIL 2012						
WEEK	M	T	W	T	F	S
13/18	30					
14	1	2	3	4	5	6
15	7	8	9	10	11	12
16	13	14	15	16	17	18
17	19	20	21	22	23	24
18	25	26	27	28	29	

Thus by this we can show that

9 # There must always exist at least one
optimal deterministic policy.



2012 APRIL

02

Monday

MAY 2012						
WK	M	T	W	T	F	S
	1	2	3	4	5	6
18	7	8	9	10	11	12
19						13
20	14	15	16	17	18	19
21	22	23	24	25	26	27
22	28	29	30	31		

DAY 093-273 WK 14

MEETINGS / APPOINTMENTS

Optimal Value Functions

→ Bellman optimality equations.

$$\pi_1 \geq \pi_2 \text{ iff } v_{\pi_1}(s) \geq v_{\pi_2}(s) \quad \forall s \in S$$

$$v_{\pi_*}(s) := \mathbb{E}[G_t | S_t = s] = \max_{\pi} v_{\pi}(s) \quad \forall s \in S$$

All optimal policies have the same } — ~~v*~~ v_*
optimal-state value function }

$$q_{\pi_*}(s, a) = \max_{\pi} q_{\pi}(s, a) \quad \forall s \in S \text{ & } a \in A$$

↳ max. action value function } — q_*

Bellman eqn for state-value func.

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

By substituting the optimal policy π_* :-

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

IMPORTANT

03

Tuesday

WK 14 DAY 094-272

APRIL

WK	M	T	W	T	F	S	S
13/18	30						
14	2	3	4	5	6	7	1
15	9	10	11	12	13	14	8
16	16	17	18	19	20	21	15
17	23	24	25	26	27	28	23

MEETINGS / APPOINTMENTS

Now π_γ will choose the action that achieves the highest value, with probability 1 and will assign probability 0 to all other actions.

$\therefore \sum_a \pi^\star(a|s)$ can be replaced by \max_a

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

\uparrow Bellman optimality eqn for V_γ

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

\uparrow Bellman optimality eqn for q_γ

Note :- for bellman eqns \Rightarrow linear system solver

will yield results

But bellman optimality eqns \Rightarrow "max" function is involved

\therefore linear system solver doesn't work.

IMPORTANT

2012 APRIL

04

Wednesday

MAY	2012					
WK	M	T	W	T	F	S
18		1	2	3	4	5
19	6	7	8	9	10	11
20	12	13	14	15	16	17
21	18	19	20	21	22	23
22	25	26	27	28	29	30
	31					

DAY 095-271 WK 14

MEETINGS / APPOINTMENTS

Bellman eqn :- $\pi, p, \gamma \rightarrow v(s)$ Bellman optimality eqn :- $\pi^*, p, \gamma \xrightarrow{?} v^*(s)$ ↳ we don't know π^*

∴ won't work.

∴ $p, \gamma \rightarrow v^*$ But our final goal is to find π^* not $v^*(s)$

↳ the optimal policy

value of the
optimal policy.

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

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

↳ particular action a that yields
the maximum value.

If there are multiple maximizing actions

we could define a stochastic optimal policy

that chooses between each of them with
some probability

MAY

JUN

IMPORTANT

APRIL 2012

05

Thursday

WK 14 DAY 096-270

APRIL

WK	M	T	W	T	F	S	S
13/18	30						
14	2	3	4	5	6	7	1
15	9	10	11	12	13	14	15
16	16	17	18	19	20	21	22
17	23	24	25	26	27	28	29

MEETINGS / APPOINTMENTS

$$\pi_\theta(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_\theta(s')]$$

10

$$\pi_\theta(s) = \underset{a}{\operatorname{argmax}} q_\theta(s, a) \rightarrow \text{easier.}$$

12

1 - Policy Evaluation & Control.

MODULE 04

↓
determining the value
func. for a specific policy

↓
finding a policy to obtain
as much reward as possible

finding a policy that
maximizes the value function.

Dynamic Programming Algorithms.

→ Use the Bellman eqn to define iterative
algorithms for both policy evaluation & control.

Policy Evaluation : $\pi \rightarrow v_\pi$

Recall : Bellman eqn gives a linear eqn for each v_π

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

IMPORTANT

2012 APRIL

06

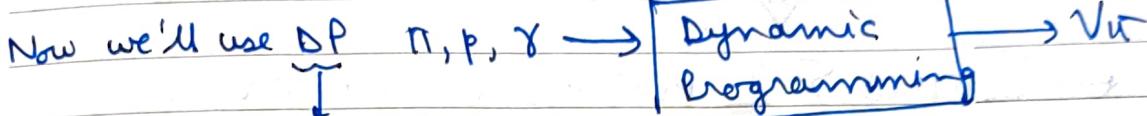
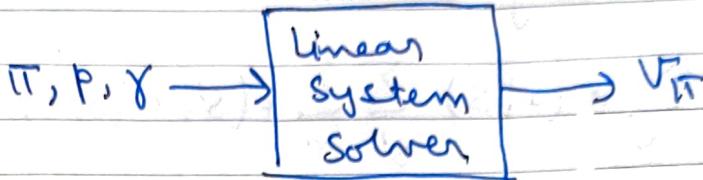
Friday

MAY 2012						
WEEK	M	T	W	T	F	S
18		1	2	3	4	5
19	6	7	8	9	10	11
20	12	13	14	15	16	17
21	18	19	20	21	22	23
22	24	25	26	27	28	29
	30	31				

DAY 097-269 WK 14

MEETINGS / APPOINTMENTS

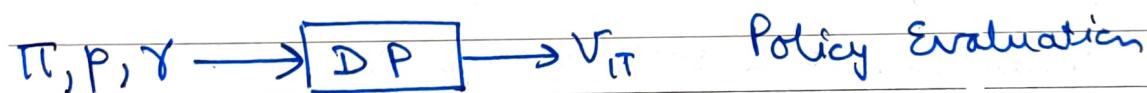
We used a linear system solver, earlier.



Iterative solution methods.

Control is the task of improving a policy. A policy π_2 is better than a policy π_1 , if π_2 is at least equal to π_1 and greater than π_1 in at least 1 state, i.e., $V_{\pi_2} \geq V_{\pi_1} \Rightarrow \pi_2$ better than π_1 .

When we can improve a policy any more, it means our current policy is the optimal policy π_\star .



DP needs access to the dynamics function p.

IMPORTANT

07

Saturday

WK 14 DAY 098-268

APRIL 2012

APRIL

WK	M	T	W	T	F	S	2012
13/18	30						
14	2	3	4	5	6	7	8
15	9	10	11	12	13	14	15
16	16	17	18	19	20	21	22
17	23	24	25	26	27	28	29

Iterative Policy Evaluation.

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



$$V_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma V_k(s')]$$

iteratively refine our estimate of the value func.
This will produce a sequence of better & better approximations to the value func.

If we reach $V_{k+1} \approx V_k$, then this is the V_{π}

for any V_0 $\underset{k \rightarrow \infty}{V_k} \rightarrow V_{\pi}$

→ Array of size equal to no. of states.

08

Sunday

WK 14 DAY 099-267

Maintain ~~to~~ 2 such arrays V , V'

\downarrow \uparrow
 V_k V_{k+1}

Input π , the policy to be evaluated.
 $V \leftarrow \vec{0}$, $V' \leftarrow \vec{0}$

loop :

$\Delta \leftarrow 0$

loop for each $s \in S$

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

$$\Delta \leftarrow \max(\Delta, |V'(s) - V(s)|)$$

$V \leftarrow V'$

until $\Delta < \theta$ (a small no.)

Output $V \approx v_{\pi}$

IMPORTANT

2012 APRIL

09

Monday

MAY 2012						
WEEK	M	T	W	T	F	S
18	1	2	3	4	5	6
19	7	8	9	10	11	12
20	14	15	16	17	18	19
21	21	22	23	24	25	26
22	28	29	30	31		

DAY 100-266 WK 15

MEETINGS / APPOINTMENTS

Policy Improvement.

$$\pi_{\text{greedy}}(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

↑ Greedy action

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

$$q_{\pi}(s, \pi'(s)) \geq q_{\pi}(s, \pi(s)) \quad \forall s \in S \Rightarrow \pi' > \pi$$

If we take action
a, according to π'
and then follows
policy π .

We take the
action under
 π

π' is better
than π

Policy Improvement Theorem

$$q_{\pi}(s, \pi'(s)) \geq q_{\pi}(s, \pi(s)) \quad \forall s \in S \rightarrow \pi' > \pi$$

$$q_{\pi}(s, \pi'(s)) > q_{\pi}(s, \pi(s)) \text{ for at least one } s \in S \rightarrow \pi' > \pi$$

The grendified policy is a strict improvement
unless the original policy was already optimal.

IMPORTANT

APRIL 2012

10

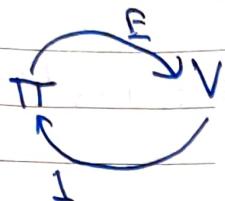
Tuesday

WK 15 DAY 101-265

MEETINGS / APPOINTMENTS

APRIL

WK	M	T	W	T	F	S	S
13/18	30						
14	1	2	3	4	5	6	7
15	9	10	11	12	13	14	15
16	16	17	18	19	20	21	22
17	23	24	25	26	27	28	29



E: Evaluation

I: improvement

10 Policy iteration algo to find the optimal policy

11 E : Evaluation (iterative policy evaluation)

I : Improvement (greedy policy improvement)

12

$$\pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} V_{\pi_2} \xrightarrow{I} \pi_3 \dots \xrightarrow{I} \pi_\star$$

 $\pi_2 > \pi_1$ $\pi_3 > \pi_2$

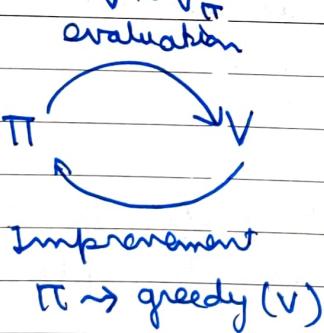
2

Sequence of better & better policy, unless the policy was already optimal.

3

Note: $\pi_2, \pi_3, \dots, \pi_\star \} \rightarrow$ Finite no. of deterministic policies.

4

If $\pi_{k+1} = \pi_k$ then $\pi_k = \pi_\star$: optimal policy

policy is greedy w.r.t. its own value function.

We alternate between evaluating the current policy and greedifying to improve the policy

IMPORTANT

2012 APRIL

2012						
MAY	M	T	W	T	F	S
WK	1	2	3	4	5	6
18						
19	7	8	9	10	11	12
20	14	15	16	17	18	19
21	21	22	23	24	25	26
22	28	29	30	31		

Wednesday

DAY 102-264 WK 15

MEETINGS / APPOINTMENTS

Pseudocode.

1. Initialization.

$V(s) \in \mathbb{R}$ and $\pi(s) \in A(s)$ arbitrarily $\forall s \in S$

2. Policy Evaluation.

Loop :

$$\Delta \leftarrow 0$$

loop for each $s \in S$:

$$v \leftarrow V(s)$$

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

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (small the no. determining accuracy of estimation)

3. Policy Improvement.

policy-stable \leftarrow true

for each $s \in S$:

$$\text{old-action} \leftarrow \pi(s)$$

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

If $\text{old-action} \neq \pi(s)$, then policy-stable \leftarrow false

If policy-stable, then stop & return

return V^* and $\pi \approx \pi^*$;

else go to 2.

IMPORTANT

12

Thursday

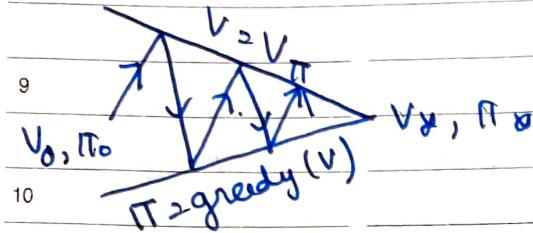
WK 15 DAY 103-263

APRIL

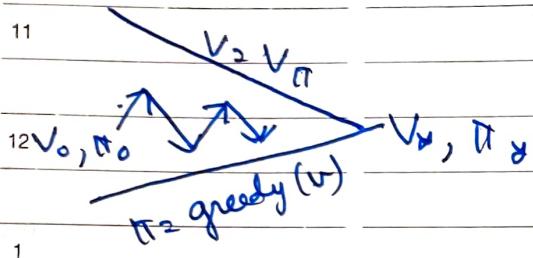
WK	M	T	W	F	S	S	2012
13/18	30						
14	2	3	4	5	6	7	1
15	9	10	11	12	13	14	15
16	16	17	18	19	20	21	16
17	23	24	25	26	27	28	27

MEETINGS / APPOINTMENTS

3. Flexibility of Policy Iteration.



The policy iteration algo runs each step all the way to completion.



Each evaluation step, brings our estimate a little closer to the value of the current policy, but not all the way.

Each policy improvement step, makes our policy a little more greedy, but not totally greedy.

3. Generalized Policy Iteration algorithms

4) 1) Value Iteration.

5. → We sweep over all the states & greedify w.r.t. the current value function.
6. But, we do not run policy evaluation to completion.

Systematic sweep } through the states } Synchronous DP.

Update the value of } states in any order } - Asynchronous DP

IMPORTANT

2012 APRIL

13

Friday

MAY	M	T	W	T	F	S	S
18	1	2	3	4	5	6	
19	7	8	9	10	11	12	13
20	14	15	16	17	18	19	20
21	21	22	23	24	25	26	27
22	28	29	30	31			

2012

DAY 104-262 WK 15

MEETINGS / APPOINTMENTS

Efficiency of DP

Monte Carlo Sampling alternative method for learning a value function.

↳ Iterative Policy evaluation is the DP solⁿ to prediction of policy evaluation problem.

↳ Monte Carlo: sample based alternative.

The value of each state can be treated as a totally independent estimation problem.

Value: $v_{\pi}(s) := \mathbb{E}_{\pi}[G_t | S_t = s]$ } Expected return from a given state.

Monte Carlo method: [→ Gather a large no. of returns under π
 → Take their average.
 → This will eventually converge to state value.]

This could deal with a lot of randomness.

Bootstrapping: we can use the other value estimates, we have already worked so hard to compute. Using the value estimates of successor states to improve our current value estimate.

↳ Much more efficient than Monte Carlo method
 estimates each value independently. IMPORTANT

APRIL 2012

14

Saturday

WK 15 DAY 105-261

APRIL

WK	M	T	W	T	F	S	S
13/18	30						
14	2	3	4	5	6	7	1
15	9	10	11	12	13	14	15
16	16	17	18	19	20	21	22
17	23	24	25	26	27	28	29

Brute force search as an alternative to Policy iteration to find optimal policy.

- Evaluates every possible deterministic policy, one at a time.
- Then pick the one with highest value.

There are finite no. of deterministic policy & there always exists an optimal deterministic policy.

Total no. of deterministic policies is exponential in the no. of states : $|A|^{|S|}$
 ↳ Brute force search can take a very long time

Policy Iteration takes polynomial time in

15 Sunday $|S|$ & $|A|$

WK 15 DAY 106-260

Therefore,

- 1) DP using bootstrapping requires much less computation than Monte Carlo
- 2) DP for policy iteration is exponentially faster than the brute force search.

DP is way better & efficient.

IMPORTANT

2012 APRIL

16

WEEK	M	T	W	T	F	S	S
18	1	2	3	4	5	6	
19	7	8	9	10	11	12	13
20	14	15	16	17	18	19	20
21	21	22	23	24	25	26	27
22	28	29	30	31			

2012

Monday

DAY 107-259 WK 16

MEETINGS / APPOINTMENTS

ADP: Approximate Dynamic Programming.

9

Dynamic Programming - Summary.

10

1. Policy Evaluation (Prediction)

$$V_{k+1}(s) := \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_k(s')]$$

11

12

2. Policy Improvement Theorem

1

3. Policy Iteration

↳ Initialization



2

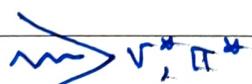
↳ Policy Evaluation : $V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_k(s')]$

3

↳ Policy Improvement : $\pi(s) = \operatorname{argmax}_a \sum_{s', r} p(s', r|s, a) [r + \gamma V_\pi(s')]$

4

4. Value Iteration



5

$$V_{k+1}(s) = \max_a \sum_{s', r} p(s', r|s, a) [r + \gamma V(s')]$$

6

MAY

JUN

IMPORTANT