

THINK OF AN ALGORITHM HAVING TWO STEPS

- 1. A POLICY IMPROVEMENT STEP: 03.
 - $\hat{\pi}(x) \in ARGMAX \quad r(x, \alpha) + \beta \left[E_{x, \alpha} \left[R(\hat{x}, \pi) \right] \right]$
- 2. A POLICY EVALUATION STEP: z_{t} . FIND $R(x,\pi) := I \sum_{x} \sum_{t=0}^{\infty} \beta^{t} r(X_{t},\pi(X_{t}))$

WE COVER TWO ALGORITHMS

1. VALUE ITERATION

2. POLICY ITERATION

VALUE MERATION TAKE V(a) =0 Hou

UNTIL CONVERGENCE DO:

 $V_{t+1}(x) = MAX \quad r(x,a) + \beta H_{x,a}[V(s^2)]$

e cod

```
def Value Iteration (V, P, r, discount):
    ''' Value Iteration – a numerical solution to a MDP
    # Arguments:
        P - P[a][x][y] gives probablity of x -> y for action a
        r - r[a][x][y] gives reward for x \rightarrow y for action a
        V - V[x] gives value for state x
        discount - a float. discount fac tor
    # Returns:
        Value function and policy from **one** value iteration
    , , ,
    number of actions = len(P)
    number of states = len(P[0])
    Q = np.zeros((number_of_actions, number_of_states))
    for in range(time):
        for a in range(number_of_actions):
            for x in range(number of states):
                Q[a][x] = np.dot(P[a][x], r[a][x]+discount*V)
        V \text{ new} = np.amax(Q, axis=0)
    pi = np.argmax(Q, axis=0)
                                                                Ħ
    return V new, pi
```

ESULT

Thrm 59. For positive programming, i.e. where all rewards are positive and the discount factor β belongs to the interval (0,1], then

 $0 \leq V_s(x) \leq V_{s+1}(x) \nearrow V(x), \quad as \quad s \to \infty.$

Here V(x) *is the optimal value function.*

PROOF IS SAME AS ne PROGRAMMINE.



Def 60 (Policy Iteration). *Given the stationary policy* Π *, we may define a new (improved) stationary policy,* $I\Pi$ *, by choosing for each* x *the action* $I\Pi(x)$ *that solves the following maximization*

$$\mathcal{I}\Pi(x) \in \underset{a \in \mathcal{A}}{\operatorname{argmax}} r(x, a) + \beta \mathbb{E}_{x, a} \left[R(\hat{X}, \Pi) \right]$$

UNTIL CONVERGENCE DO

$$||_{n+1} = ||_{n+1}$$

FINDING $R(2,\pi)$: WE KNOW FROM MARKOJ CHIAINS (LECTURE 2) THAT

$$R(x) = r(x) + \beta E_{xc} R(\hat{x})$$
HERE
$$R(x) = R(x,\pi), \quad r(x) = r(x,\pi), \quad P_{xg} = P_{xg}^{\pi(a)}.$$
INSTETRET AS VECTORS & MATTRICES.

$$\frac{R}{R} = \underline{r} + \beta PR \quad (\Longrightarrow) \quad \frac{R}{R} = (\underline{r} - \beta R)^{T} \underline{r}$$

$$\downarrow S \quad JUST \quad MATRIX \quad ALGEBRA$$

ome Code:

def Policy_Iteration(pi,P,r,discount):
 ''' Policy Iteration - a numerical solution to a MDP

Arguments:

P - P[a][x][y] gives probablity of x -> y for action a r - r[a][x][y] gives reward for x -> y for action a pi - pi[x] gives action for state x discount - disount factor

Returns:

. . .

policy from **one** policy iteration
value function of input policy

Collate array of states and actions

number_of_actions, number_of_states = len(P), len(P[0])
Actions, States = np.arange(number_of_actions), np.arange(
number_of_states)

```
# Get transitions and rewards of policy pi
P_pi = np.array([P[pi[x]][x] for x in States ])
r_pi = np.array([r[pi[x]][x] for x in States])
Er_pi = [ np.dot(P_pi[x], r_pi[x]) for x in States]
```

```
# Calculate Value of pi
I = np.identity(number_of_states)
A = I - discount * P_pi
R_pi = np.linalg.solve(A, Er_pi)
```

```
# Calculate Q_factors of pi
Q = np.zeros((number_of_actions,number_of_states))
for a in range(number_of_actions):
    for x in range(number_of_states):
        Q[a][x] = np.dot(P[a][x],r[a][x]+discount*R_pi)
```

```
# policy iteration update
pi_new = np.argmax(Q, axis=0)
```

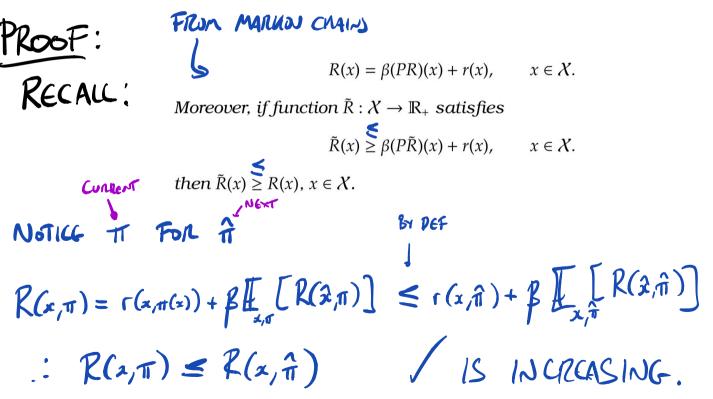
return pi_new, R_pi

Thrm 61. Under Policy Iteration

 $R(x, \Pi_{n+1}) \ge R(x, \Pi_n)$

and, for bounded programming,

```
R(x, \Pi_n) \nearrow V(x) as n \to \infty
```



To PRove CONVERCENCE TO OPTIMAL POLICY
LET
$$M_{t} = \sum_{s=0}^{t-1} \beta^{s} r(x, \pi^{*}(x)) + \beta^{t} R(x, \pi_{T-t})$$

$$M_{t} = \sum_{s=0}^{t-1} \beta^{s} r(x, \pi^{*}(x)) + \beta^{t} R(x, \pi_{T-t})$$

$$The Do TT_{T-t}$$

$$\mathbb{E}\left[M_{t+1} - M_{t} \right]_{0}^{o} = \beta^{t} \mathbb{E}\left[R(x_{t+1}, \pi_{T-t-1}) + r(x_{t}, \pi^{*}) - R(x, \pi_{T-t})\right]_{0}^{1}$$

$$\mathbb{E}\left[M_{t+1} - M_{t} \right]_{0}^{o} = \beta^{t} \mathbb{E}\left[R(x_{t+1}, \pi_{T-t-1}) + r(x_{t}, \pi^{*}) - R(x, \pi_{T-t})\right]_{0}^{1}$$

$$\mathbb{E}\left[M_{t+1} - M_{t} \right]_{0}^{o} = \beta^{t} \mathbb{E}\left[R(x_{t+1}, \pi_{T-t-1}) + r(x_{t}, \pi^{*}) - R(x, \pi_{T-t})\right]_{0}^{1}$$

.. ME IS SUPER MJ.

 $\mathbb{R}(x,\pi_{\varepsilon}) \neq \mathbb{E}[M_{0}] \ge \mathbb{E}[M_{\tau}]$ $\cong \mathbb{R}_{\tau}(x,\pi^{*}) \longrightarrow \mathbb{R}(x,\pi^{*}) \longrightarrow \mathbb{R}(x,\pi^{*}) \longrightarrow \mathbb{R}(x,\pi^{*}) \longrightarrow \mathbb{R}(x,\pi^{*})$