## **Recitation 3: Homework 1**

## Yafei Hu and Justin Kiefel

#### 1.1: Contraction Mapping

An operator F on a normed vector space  $\mathscr{X}$  is a  $\gamma$ -contraction, for  $0 < \gamma < 1$  provided for all  $x, y \in \mathscr{X}$ :

 $||F(x) - F(y)|| \le \gamma ||x - y||$ 

Theorem (Contraction mapping) For a  $\gamma$ -contraction F in a complete normed vector space  $\mathcal{X}$ :

- F converges to a unique fixed point in  $\mathcal{X}$ ,
- at a linear convergence rate  $\gamma$ .

Fixed point definition:

A point/vector  $x \in \mathcal{X}$  is a fixed point of an operator F if F(x) = x

Banach Fixed Point Theorem:

If F is a  $\ \gamma$  -contraction mapping, then:

- F has a **unique** fixed point
- $\forall x_0 \in \mathcal{X}$ , the sequence  $x_{n+1} = F(x_n)$  converges to  $x^*$  in a geometric fashion:

$$||x_n - x^*|| \le \gamma^n ||x_0 - x^*||$$

thus 
$$\lim_{n \to \infty} ||x_n - x^*|| \le \lim_{n \to \infty} (\gamma^n ||x_0 - x^*||) = 0$$

#### **Policy Evaluation**

```
Iterative Policy Evaluation, for estimating V \approx v_{\pi}
```

```
Input \pi, the policy to be evaluated

Algorithm parameter: a small threshold \theta > 0 determining accuracy of estimation

Initialize V(s), for all s \in S^+, arbitrarily except that V(terminal) = 0

Loop:

\Delta \leftarrow 0

Loop for each s \in S:

v \leftarrow V(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 - V(s)|)

until \Delta < \theta
```

#### **Policy Iteration**

Note the

difference

between sync

and async Policy

Iteration

Policy Iteration (using iterative policy evaluation) for estimating  $\pi \approx \pi_*$ 1. Initialization  $V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in S$ 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  $\Delta < \theta$  (a small positive number determining the accuracy of estimation) 3. Policy Improvement policy-stable  $\leftarrow true$ For each  $s \in S$ : old-action  $\leftarrow \pi(s)$  $\pi(s) \leftarrow \operatorname{arg\,max}_{a} \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$ If old-action  $\neq \pi(s)$ , then policy-stable  $\leftarrow$  false If *policy-stable*, then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

#### Value Iteration

	Value Iteration, for estimating $\pi \approx \pi_*$
Note the difference between sync and async Value Iteration	Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation Initialize $V(s)$ , for all $s \in S^+$ , arbitrarily except that $V(terminal) = 0$ Loop: $\Delta \leftarrow 0$ Loop for each $s \in S$ : $v \leftarrow V(s)$ $V(s) \leftarrow \max_a \sum_{s',r} p(s', r   s, a) [r + \gamma V(s')]$ $\Delta \leftarrow \max(\Delta,  v - V(s) )$ until $\Delta < \theta$ Output a deterministic policy, $\pi \approx \pi_*$ , such that $\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s', r   s, a) [r + \gamma V(s')]$

Synchronous and Asynchronous Policy Iteration/Value Iteration

- Synchronous value iteration stores two copies of value function
  - for all s in  $\mathcal S$

$$\mathbf{v}_{new}(s) \leftarrow \max_{a \in \mathcal{A}} \left( r(s, a) + \gamma \sum_{s' \in \mathcal{D}} p\left(s' \mid s, a\right) \mathbf{v}_{old}(s') \right)$$

 $v_{old} \leftarrow v_{new}$ 

- · In-place value iteration only stores one copy of value function
  - for all s in S

$$\boldsymbol{\nu(s)} \leftarrow \max_{a \in \mathcal{A}} \left( r(s, a) + \gamma \sum_{s' \in \mathcal{S}} T(s' | s, a) \boldsymbol{\nu(s')} \right)$$

#### Synchronous and Asynchronous Policy Iteration/Value Iteration

Example: synchronous and randperm asynchronous Value Iteration in 4 × 4 Frozen Lake Environment

#### Synchronous

Asynchronous

Previous value function

Updated value function

Updated value function

teration	1	

0.	0.	0.	0.
0.	0.	0.	0.
0.	0.	0.	0.
0.	0.	0.	0.

0.	0.	0.	0.
1.	0.	0.	0.
0.	0.	0.	0.
0.	0.	0.	0.

0.	0.	0.
0.	0.	0.
0.	0.	0.
0.	0.	0.
	0. 0.	0.     0.       0.     0.

0.	0.	0.	0.
1.	0.	0.	0.
0.	0.	0.	0.
0.	0.	0.	0.

0.9	0.	0.	0.
1.	0.	0.	0.
0.9	0.	0.	0.
0.	0.	0.	0.

0.9	0.	0.	0.
1.	0.	0.	0.
0.9	0.	0.	0.
0.81	0.729	0.	0.

Iteration	3
-----------	---

Iteration 2

0.9	0.	0.	0.	
1.	0.	0.	0.	
0.9	0.	0.	0.	
0.	0.	0.	0.	

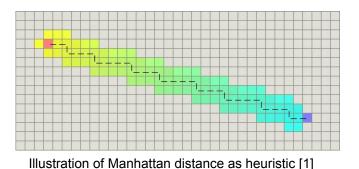
0.9	0.	0.	0.
1.	0.	0.	0.
0.9	0.	0.	0.
0.81	0.	0.	0.

0.9	0.	0.	0.
1.	0.	0.	0.
0.9	0.	0.	0.
0.81	0.729	0.6561	0.

#### Problem 1.5: Manhattan distance as heuristic function

- 1. Compute the heuristic function for all the states in advance
- Sweep the states ordered by the pre-computed heuristic function (in this case, Manhattan distance)
   Manhattan distance between some state *S* and goal *G*

$$d(S,G) = \sum_{i=1}^{n} |S_i - G_i|$$
, where  $n = 2$ 



More about heuristic: [1] http://theory.stanford.edu/~amitp/GameProgramming/Heuristics.html

### Problem 1 Q & A

1. Will I get the same iteration number in PI/VI and same policy every time?

Yes, you are expected to get same results every time you run the experiment.

# **Problem 2: Bandits**

### **Estimating Expected Reward**

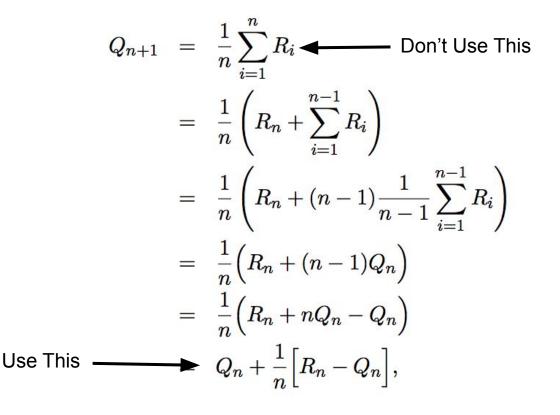
$$\mathbb{E}\{R_t\} = \frac{1}{20} \Sigma_{k=1}^{20} R_t^k$$

- Average of rewards received at a given time step
- Unbiased
- High Variance

$$\mathbb{E}\{R_t\} = \frac{1}{20} \Sigma_{k=1}^{20} \mathbb{E}\{r^k (A_t^k) | \pi_t^k\}$$

- Average of expected rewards conditioned on the policy
- Unbiased
- Lower Variance
- Remember to still use  $R_t$  for the agent's update

#### **Efficient Q-Updates**



#### Problem 2.7 - Correlated Rewards

I.I.D. Rewards

**Correlated Rewards** 

$$r(k) \sim \mathcal{N}(\mu, \sigma^2) \ \forall k \in [K]$$

 $[r(1) \dots r(K)]^T \sim \mathcal{N}(\mu_0, \Sigma_0)$ Non-Diagonal

# Questions?