

## 1 Problem3 plots

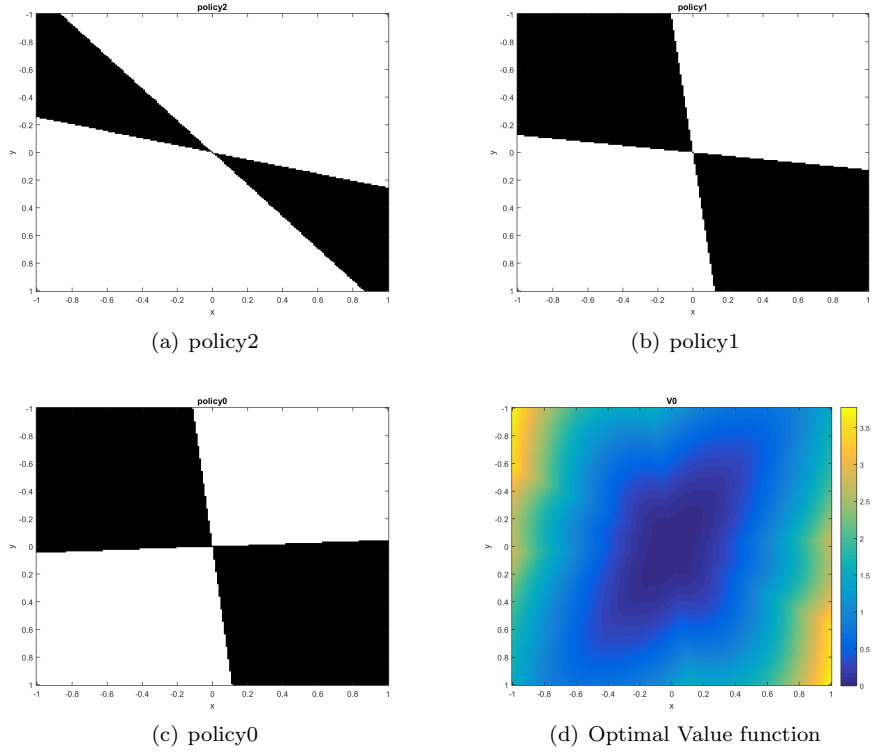


Figure 1: All the plots for problem 3b. Note that the dark region in all policies means optimal choice is select policy 1 and while region means select 2.

## 2 Problem 4

### 2.1 Problem Formulation

Due to the assumption of no cycles with negative cost, the optimal path need not have more than  $|V|$  elements. Then we can formulate the DSP problem as a DFS with  $T := |V| - 1$  stages:

- 1.State space:  $\mathcal{X}_0 := \{s\}, \mathcal{X}_T := \{\tau\}, \mathcal{X}_t := \mathcal{V} \setminus \{\tau\}$  for  $t = 1, \dots, T - 1$
- 2.Control space:  $u_{T-1} := \{\tau\}$  and  $u_t := \mathcal{V} \setminus \{\tau\}$  for  $t = 0, \dots, T - 2$
- 3.Dynamics(motion model):  $x_{t+1} = u_t$  for  $u_t \in \mathcal{U}_t, t = 0, \dots, T - 1$
- 4.Costs:  $q(\tau) := 0$  and  $\ell_t(x_t, u_t) = c_{x_t, u_t}$  for  $t = 0, \dots, T - 1$

### 2.2 Technical Approach

By implementing Dynamic Programming algorithm on DSP problem, we can have the follow algorithm:

Because the DSP problem is symmetric: an optimal path from s to t is also a shortest path from t to s, where all arc directions are flipped.

For simplicity, we only need to implement label correcting (LC) algorithm. It is a general algorithm for SP problems that does not necessarily visit every node of the graph. It is shown as follow:

---

**Algorithm 1** Deterministic Shortest Path via Dynamic Programming

---

```
1: Input: node set  $\mathcal{V}$ , start  $s \in \mathcal{V}$ , goal  $\tau \in \mathcal{V}$ , and costs  $c_{ij}$  for  $i, j \in \mathcal{V}$ 
2:  $T = |\mathcal{V}| - 1$ 
3:  $V_T(\tau) = 0$ 
4:  $V_{T-1}(i) = c_{i,\tau}, \quad \forall i \in \mathcal{V} \setminus \{\tau\}$ 
5:  $\pi_{T-1}(i) = \tau, \quad \forall i \in \mathcal{V} \setminus \{\tau\}$ 
6: for  $t = (T-2), \dots, 0$  do
7:    $Q_t(i, j) = c_{ij} + V_{t+1}(j), \quad \forall i, j \in \mathcal{V} \setminus \{\tau\}$ 
8:    $V_t(i) = \min_{j \in \mathcal{V} \setminus \{\tau\}} Q_t(i, j), \quad \forall i \in \mathcal{V} \setminus \{\tau\}$ 
9:    $\pi_t(i) = \arg \min_{j \in \mathcal{V} \setminus \{\tau\}} Q_t(i, j), \quad \forall i \in \mathcal{V} \setminus \{\tau\}$ 
10:  if  $V_t(i) = V_{t+1}(i), \forall i \in \mathcal{V} \setminus \{\tau\}$  then
11:    break
```

---

Figure 2: Dynamic Programming algorithm for DSP problem

---

**Algorithm 2** Deterministic Shortest Path via Forward Dynamic Programming

---

```
1: Input: node set  $\mathcal{V}$ , start  $s \in \mathcal{V}$ , goal  $\tau \in \mathcal{V}$ , and costs  $c_{ij}$  for  $i, j \in \mathcal{V}$ 
2:  $T = |\mathcal{V}| - 1$ 
3:  $V_0^F(s) = 0$ 
4:  $V_1^F(j) = c_{sj}, \quad \forall j \in \mathcal{V} \setminus \{s\}$ 
5: for  $t = 2, \dots, T$  do
6:    $V_t^F(j) = \min_{i \in \mathcal{V} \setminus \{s\}} (c_{ij} + V_{t-1}^F(i)), \quad \forall j \in \mathcal{V} \setminus \{s\}$ 
7:   if  $V_t^F(\tau) < \infty$  then
8:     break
```

---

Figure 3: Forward Dynamic Programming algorithm for DSP problem

---

**Algorithm 3** Label Correcting Algorithm

---

```
1:  $\text{OPEN} \leftarrow \{s\}, g_s = 0, g_i = \infty$  for all  $i \in \mathcal{V} \setminus \{s\}$ 
2: while  $\text{OPEN}$  is not empty do
3:   Remove  $i$  from  $\text{OPEN}$ 
4:   for  $j \in \text{Children}(i)$  do
5:     if  $(g_i + c_{ij}) < g_j$  and  $(g_i + c_{ij}) < g_\tau$  then
6:        $g_j \leftarrow (g_i + c_{ij})$ 
7:        $\text{Parent}(j) \leftarrow i$ 
8:       if  $j \neq \tau$  then
9:          $\text{OPEN} \leftarrow \text{OPEN} \cup \{j\}$ 
```

---

Figure 4: Label correcting algorithm for DSP problem

## 2.3 Results

### 2.3.1 p1

minimum cost path: 42 43 44 53 61 70 79 80 81 82 83 84 85 86 87 98 109

optimal cost-to-go values: 16.00 15.00 14.00 13.00 12.00 11.00 10.00 9.00 8.00 7.00 6.00 5.00 4.00 3.00 2.00 1.00 0.00

### 2.3.2 p2

minimum cost path: 20 31 42 53 64 75 86 97 108 109 120

optimal cost-to-go values: 67.66 67.41 66.94 66.11 64.53 61.38 55.53 46.00 32.62 16.53 0.00

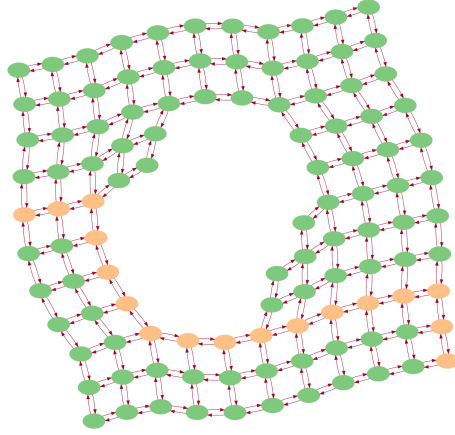


Figure 5: Plot for problem 1.

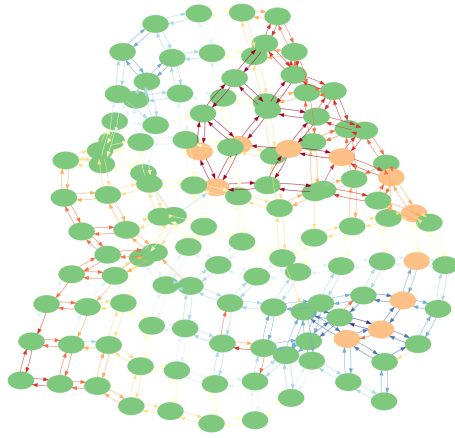


Figure 6: Plot for problem 2.

### 2.3.3 p3

minimum cost path: 2 5 9 13 18 23

optimal cost-to-go values: 6.24 4.83 3.41 2.41 1.41 0.00

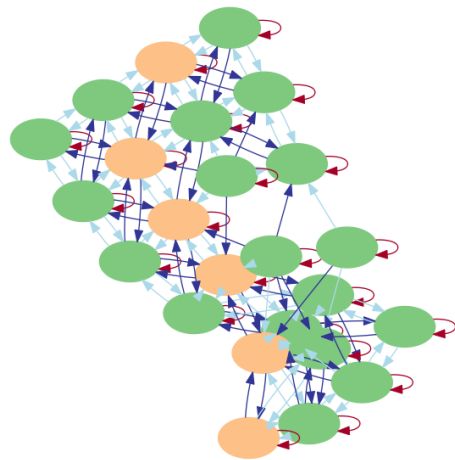


Figure 7: Plot for problem 3.

#### 2.3.4 p4

minimum cost path: 11 22 33 44 55 66 67 68 69 70 71 72 83 84 85 96 97 108 109 120

optimal cost-to-go values: 82.52 82.27 81.79 80.97 79.61 77.52 74.48 70.33 64.93 58.18 50.12 41.16 32.30  
24.61 18.83 14.55 10.69 6.89 3.21 0.00

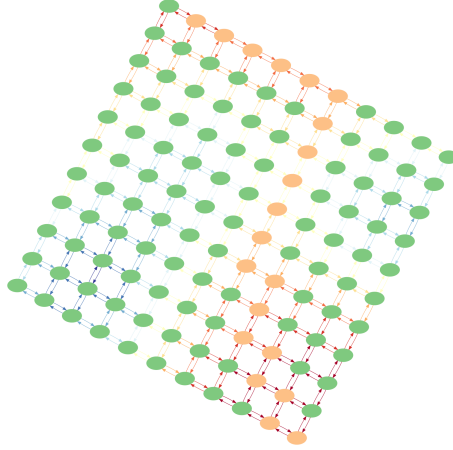


Figure 8: Plot for problem 4.

#### 2.3.5 p5

minimum cost path: 0 11 22 23 24 25 26 27 28 29 30 31 32 41 50 58 67 76 87 98 109

optimal cost-to-go values: 41.66 39.66 37.66 35.66 33.66 31.66 29.66 27.66 25.66 23.66 21.66 19.66 17.66  
15.66 13.66 11.66 9.66 7.66 5.66 2.45 0.00

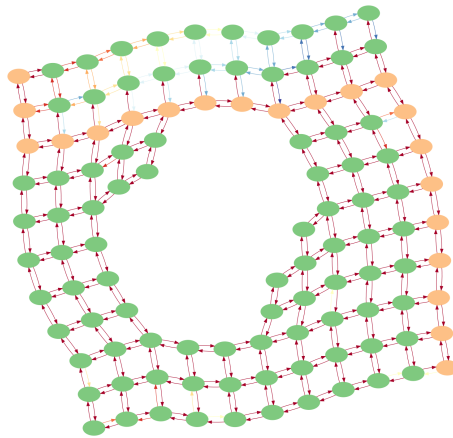


Figure 9: Plot for problem 5.

#### 2.3.6 p6

minimum cost path: 27 38 37 45 46 109 120 119 118

optimal cost-to-go values: 43.18 42.49 42.21 41.96 29.97 22.60 16.67 4.04 0.00

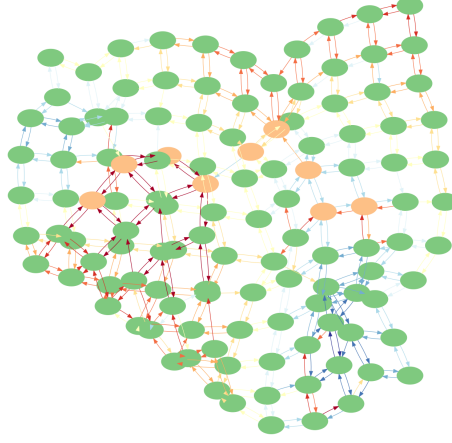


Figure 10: Plot for problem 6.

## 3 Problem 5

### 3.1 Problem Formulation

#### 3.1.1 POMDP

We first define a POMDP  $(\mathcal{X}, \mathcal{U}, \mathcal{Z}, p_f, p_h, \ell, \gamma)$  problem as follows:

1.State space  $\mathcal{X}$ : We keep track of the number of papers(np), scissors(ns) and rocks(nr) that opponent gives. Also, we keep track of our total score differential(Sd). Our state space are extremely huge. At time step t, we have  $(2*t+1)$  number of possible scores and  $\binom{t+2}{2}$  number of possible number of combinations of total number of (paper ,scissors rocks). Thus,

$$(2 * 0 + 1) \binom{2}{2} + (2 * 1 + 1)(2) + \dots + (2 * 100 + 1) \binom{102}{2} = 6 \binom{103}{4} + \binom{103}{3} \quad (1)$$

number of states in total.

2.Control space  $\mathcal{U}$ : Control space is just what you move at each round.(paper, scissor or rock)

3.Observation space  $\mathcal{Z}$ : Observation space is win(1), draw(0) or lose(-1) at each round.

4.Observation model  $p_h$ : After each round, we have an observation of the game, which is win(+1), draw(0) or loss(-1), which is an absolute certain value.

5.Dynamics(motion model)  $p_f$ : Given history of opponent and my movement at time t, we need to formulate the distribution of next state  $x_{t+1} \sim p_f(\cdot | x_t, u_t)$ . Note that:

$$\begin{aligned} &Pr(\text{rock} | X_r = x_r, X_p = x_p, X_s = x_s) = \\ &Pr(\text{rock} | F = \text{rock}, X_r = x_r, X_p = x_p, X_s = x_s)Pr(F = \text{rock} | X_r = x_r, X_p = x_p, X_s = x_s) \\ &+ Pr(\text{rock} | F = \text{paper}, X_r = x_r, X_p = x_p, X_s = x_s)Pr(F = \text{paper} | X_r = x_r, X_p = x_p, X_s = x_s) \\ &+ Pr(\text{rock} | F = \text{scissors}, X_r = x_r, X_p = x_p, X_s = x_s)Pr(F = \text{scissors} | X_r = x_r, X_p = x_p, X_s = x_s) \end{aligned}$$

Since

$$\Pr(\text{rock} | F = \text{rock}, X_r = x_r, X_p = x_p, X_s = x_s)_s = 0.5 \quad (2)$$

$$\Pr(F = \text{rock} | X_r = x_r, X_p = x_p, X_s = x_s) = \frac{\Pr(X_r = x_r, X_p = x_p, X_s = x_s | F = \text{rock}) * \Pr(F = \text{rock})}{\Pr(X_r = x_r, X_p = x_p, X_s = x_s)} \quad (3)$$

Assume that the prior  $\Pr(F = \text{rock}) = \Pr(F = \text{scissor}) = \Pr(F = \text{paper}) = \frac{1}{3}$ ,

Also,

$$\begin{aligned} &\Pr(F = \text{rock} | X_r = x_r, X_p = x_p, X_s = x_s) + \Pr(F = \text{scissor} | X_r = x_r, X_p = x_p, X_s = x_s) \\ &\quad + \Pr(F = \text{paper} | X_r = x_r, X_p = x_p, X_s = x_s) = 1 \end{aligned} \quad (4)$$

Here we have

$$\Pr(F = \text{rock} | X_r = x_r, X_p = x_p, X_s = x_s) = \frac{\binom{x_r + x_p + x_s}{x_r}}{\binom{x_r + x_p + x_s}{x_r} + \binom{x_r + x_p + x_s}{x_p} + \binom{x_r + x_p + x_s}{x_s}} \quad (5)$$

In this case, we have

$$\begin{aligned} &\Pr(\text{rock} | X_r = x_r, X_p = x_p, X_s = x_s) = \\ &\frac{0.5 * \binom{x_r + x_p + x_s}{x_r} + 0.25 * \binom{x_r + x_p + x_s}{x_p} + 0.25 * \binom{x_r + x_p + x_s}{x_s}}{\binom{x_r + x_p + x_s}{x_r} + \binom{x_r + x_p + x_s}{x_p} + \binom{x_r + x_p + x_s}{x_s}} \end{aligned} \quad (6)$$

We have calculate the conditional probability on State  $X_t$ , then we can formulate the motion model by updating  $X_{t+1}$ . Without the loss of generality, assume we have the movement  $u_t = \text{paper}$ , then

$$\begin{aligned} &\Pr(X_r^{t+1} = x_r + 1, X_p^{t+1} = x_p, X_s^{t+1} = x_s, Sd^{t+1} = sd + 1 | U_t = \text{paper}) = \\ &\Pr(\text{rock} | X_r^t = x_r, X_p^t = x_p, X_s^t = x_s, Sd^t = sd) = \\ &\frac{0.5 * \binom{x_r + x_p + x_s}{x_r} + 0.25 * \binom{x_r + x_p + x_s}{x_p} + 0.25 * \binom{x_r + x_p + x_s}{x_s}}{\binom{x_r + x_p + x_s}{x_r} + \binom{x_r + x_p + x_s}{x_p} + \binom{x_r + x_p + x_s}{x_s}} \end{aligned} \quad (7)$$

$$\begin{aligned} &\Pr(X_r^{t+1} = x_r, X_p^{t+1} = x_p + 1, X_s^{t+1} = x_s, Sd^{t+1} = sd | U_t = \text{paper}) = \\ &\Pr(\text{paper} | X_r^t = x_r, X_p^t = x_p, X_s^t = x_s, Sd^t = sd) = \\ &\frac{0.5 * \binom{x_r + x_p + x_s}{x_p} + 0.25 * \binom{x_r + x_p + x_s}{x_r} + 0.25 * \binom{x_r + x_p + x_s}{x_s}}{\binom{x_r + x_p + x_s}{x_r} + \binom{x_r + x_p + x_s}{x_p} + \binom{x_r + x_p + x_s}{x_s}} \end{aligned} \quad (8)$$

$$\begin{aligned} &\Pr(X_r^{t+1} = x_r, X_p^{t+1} = x_p, X_s^{t+1} = x_s + 1, Sd^{t+1} = sd - 1 | U_t = \text{paper}) = \\ &\Pr(\text{scissor} | X_r^t = x_r, X_p^t = x_p, X_s^t = x_s, Sd^t = sd) = \\ &\frac{0.5 * \binom{x_r + x_p + x_s}{x_s} + 0.25 * \binom{x_r + x_p + x_s}{x_p} + 0.25 * \binom{x_r + x_p + x_s}{x_r}}{\binom{x_r + x_p + x_s}{x_r} + \binom{x_r + x_p + x_s}{x_p} + \binom{x_r + x_p + x_s}{x_s}} \end{aligned} \quad (9)$$

which is similar for  $U_t = \text{rock}, \text{scissor}$ .

6. Initial state  $x_0: x_0 = \{x_s = x_r = x_p = 0, sd = 0\}$

7.Planning horizon:  $T = 100$ .

8.Costs: We only have the terminal cost here, we define  $q(Sd) := -Sd$ , which is the negtive of the score differential(Sd).

Here is the summation of POMDP formulation.

$$\begin{aligned}
& \min_{\pi_{0:T-1}} E \left[ \gamma^T q(x_T) + \sum_{t=0}^{T-1} \gamma^t \ell_t(x_t, u_t) \right] \\
& \text{s.t. } x_{t+1} \sim p_f(\cdot | x_t, u_t), \quad t = 0, \dots, T-1 \\
& \quad z_{t+1} \sim p_h(\cdot | x_t), \quad t = 0, \dots, T-1 \\
& \quad u_t \sim \pi_t(\cdot | i_t), \quad t = 0, \dots, T-1 \\
& \quad x_0 \sim b_0(\cdot) \equiv \text{prior pdf over the hidden state } x_0
\end{aligned} \tag{10}$$

### 3.1.2 MDP

Then, we use the equavalence to formulate  $\text{MDP}(\mathcal{B}, \mathcal{U}, p_\psi, \rho, \gamma)$ . Here we define our MDP problem as follows:

- 1.State space  $\mathcal{B} := \mathcal{P}(\mathcal{X})$  is the continuous space of pmfs over  $\mathcal{X}$ , which has been defined above.
- 2.Control space: Similar as POMDP, control space is just what you move at each round.(paper, scissor or rock)
- 3.Transition model  $p_\psi$ : In MDP, we need to eliminate our observation space. The transformed motion model is the Bayes filter  $b_{t+1} = \psi(b_t, u_t, z_t)$ , where  $z_t$  plays the role of noise or in probabilistic terms. Given the distribution of  $Z_{t+1}, b_t, u_t$ , we can easily formulate  $b_{t+1}$ . Because at each round, given my movement, the probability of win, draw and lose is the probability of +1,0,-1 to Sd, and according to  $u_t$ , we can easily formulate what our opponent's movement as follows:

Transition model	$u_t = \textit{paper}$	$u_t = \textit{scissor}$	$u_t = \textit{rock}$
$z_t = \textit{win}(+1)$	$x_r + = 1, sd + = 1$	$x_p + = 1, sd + = 1$	$x_s + = 1, sd + = 1$
$z_t = \textit{draw}(0)$	$x_p + = 1$	$x_s + = 1$	$x_r + = 1$
$z_t = \textit{loss}(-1)$	$x_s + = 1, sd - = 1$	$x_r + = 1, sd - = 1$	$x_p + = 1, sd - = 1$

4.Initial state  $b_0$ :  $b_0 = \{x_s = x_r = x_p = 0, sd = 0\}$

5.Planning horizon:  $T = 100$ .

6.Costs: We only have the terminal cost here, we define  $q(Sd) := -Sd$ , which is the negative of the score differential(Sd).

Here is the summation of MDP formulation:

$$\min_{\pi_{0:T-1}} V_0^\pi(b_0) = E \left[ \gamma^T \rho_T(b_T) + \sum_{t=0}^{T-1} \gamma^t \rho_t(b_t, u_t) \right] \tag{11}$$

$$\begin{aligned}
& \text{s.t. } b_{t+1} = \psi(b_t, u_t, z_{t+1}), \quad t = 0, \dots, T-1 \\
& \quad z_{t+1} \sim \eta(\cdot | b_t, u_t), \quad t = 0, \dots, T-1 \\
& \quad u_t \sim \pi_t(\cdot | b_t), \quad t = 0, \dots, T-1
\end{aligned} \tag{12}$$

## 3.2 Technical Approach

For simplicity, we first propose the optimal policy and then prove its optimality.

### 3.2.1 Optimal Policy

By applying dynamic programming algorithm on our model, we can formulate the optimal policy as follows: Each round, we choose the movement with largest count as opponent's favor. Then, we select our movement according to opponent's favor to beat him.(i.e.  $F^t = \operatorname{argmax}\{x_p^t, x_s^t, x_r^t\}$ ) We randomly choose movement when 3 action history are equal.

### 3.2.2 Prove of Optimality

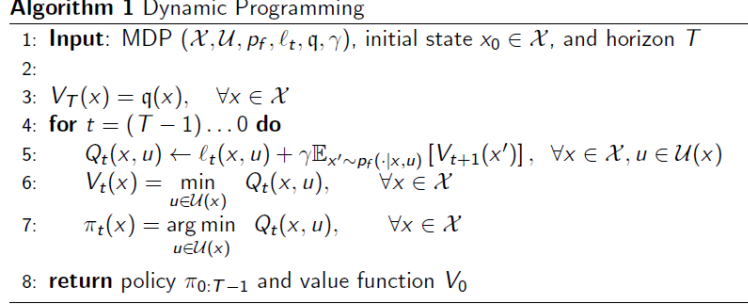


Figure 11: DP algorithm

Inferring DP algorithm, we only have terminal cost, so here  $\ell_t(x, u) = 0$ . So, our optimal policy should be chosen from

$$\pi_t(x) = \arg \min_{u \in \mathcal{U}(x)} Q_t(x, u), \quad \forall x \in \mathcal{X} \quad (13)$$

which is equal to

$$\pi_t(x) = \arg \min_{u \in \mathcal{U}(x)} E_{x' \sim p_f(\cdot|x, u)} [V_{t+1}(x')] = \arg \max_{u \in \mathcal{U}(x)} E_{x' \sim p_f(\cdot|x, u)} [Sd'], \quad \forall x \in \mathcal{X} \quad (14)$$

According to the motion model defined above,

$$\begin{aligned}
Q_t(x, u) &= \Pr(Sd^{t+1} = sd + 1 | U_t) * (sd + 1) + \Pr(Sd^{t+1} = sd | U_t) * (sd) \\
&\quad + \Pr(Sd^{t+1} = sd - 1 | U_t) * (sd - 1) \\
&= sd + \Pr(Sd^{t+1} = sd + 1 | U_t) - \Pr(Sd^{t+1} = sd - 1 | U_t)
\end{aligned} \quad (15)$$

Thus,

$$\begin{aligned}
\pi_t(x) &= \arg \max_{u \in \mathcal{U}(x)} E_{x' \sim p_f(\cdot|x, u)} [Sd'], \quad \forall x \in \mathcal{X} \\
&= \arg \max_{u \in \mathcal{U}(x)} [\Pr(Sd^{t+1} = sd + 1 | U_t) - \Pr(Sd^{t+1} = sd - 1 | U_t)] \\
&= \arg \max_{u \in \mathcal{U}(x)} \left[ \begin{pmatrix} x_r + x_p + x_s \\ x_{u_{win}} \end{pmatrix} - \begin{pmatrix} x_r + x_p + x_s \\ x_{u_{lose}} \end{pmatrix} \right]
\end{aligned} \quad (16)$$

So, we can easily observed that we need to choose a movement( $u_{win}$ ) which can win opponent's top number of action( $x_{u_{win}}$ ).



### 3.3 Result

We can observed in the plot that our optimal policy is indeed optimal. With increasing number of games, the mean score differential continuous increasing, while stochastic and deterministic policy are around 0. Also, optimal policy tend to be a straight line in mean plot when number of games are becoming large, since our evidence are accumulating, we are more absolute about our inference result. Also, the mean score differential of our optimal policy approximate 25 out of 100 games, which means 50% chance of winning, 25% drawing and 25% of losing.

From standard deviation plot, we can observe that all policy are showing increasing trend, since our state space is increasing as time increase. Moreover, the optimal policy has largest standard deviation, which might because our optimal policy is dependent of our opponent's behavior, if our opponent's movement are abnormal at very beginning, our policy might choose some bad result.

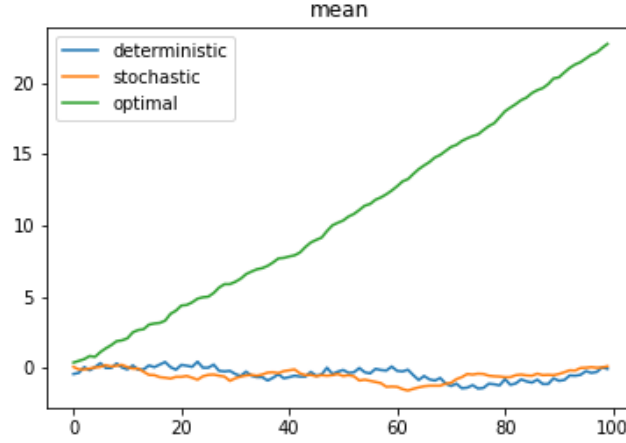


Figure 12: Mean over number of games.

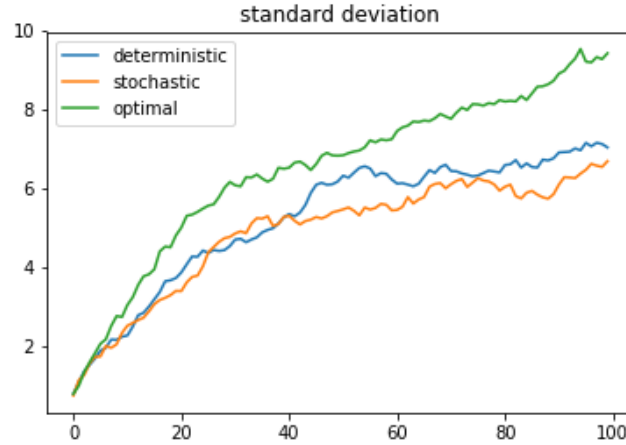


Figure 13: Standard Deviation over number of games.