

# CS 181 Machine Learning

## Practical 4 Report, Team *la Dernière Dame M*

(Jeremiah) Zhe Liu<sup>1</sup>, (Vivian) Wenwan Yang<sup>2</sup>, and Jing Wen<sup>1</sup>

<sup>1</sup>Department of Biostatistics, Harvard School of Public Health

<sup>2</sup>Department of Computational Science and Engineering, SEAS

May 6, 2015

### 1 Problem Description

Set in a *Flappy Bird*-type game *Swingy Monkey*, our current learning task is to estimate the optimal policy function  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  such that the expectation of reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$  is maximized, i.e. if define a stochastic process of game state  $\{s_t\}$  with unknown transition probability  $P(s_{t+1} | s_1, \dots, s_t, a_1, \dots, a_t)$ , we aim to identify a  $\pi^*$  such that

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left( \sum_{s \in \mathbf{p}} R(s, \pi(s)) \middle| \mathbf{p} \right)$$

where  $\mathbf{p}$  a sample path of  $\{S_t\}$ .

In current setting, the state and action spaces are defined as:

$$\begin{aligned} \mathcal{S} &= [\text{Tree}_{\text{dist}} \quad \text{Tree}_{\text{top}} \quad \text{Tree}_{\text{bot}} \quad \text{Monkey}_{\text{vel}} \quad \text{Monkey}_{\text{top}} \quad \text{Monkey}_{\text{bot}}] \subset \mathbb{R}^6 \\ \mathcal{A} &= [\text{NoJump} \quad \text{Jump}] \end{aligned}$$

Note that  $[\text{Monkey}_{\text{top}}, \text{Monkey}_{\text{bot}}, \text{Tree}_{\text{top}}, \text{Tree}_{\text{bot}}]$  are in fact bounded by screen size (600 pxls).

The reward function can be partially described as:

$$R : \begin{bmatrix} \text{pass\_tree} \\ \text{hit\_trunk} \\ \text{hit\_edge} \\ \text{otherwise} \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ -5 \\ -10 \\ 0 \end{bmatrix}$$

where  $[\text{pass\_tree}, \text{hit\_trunk}, \text{hit\_edge}]$  are unknown boolean functions of  $s \in \mathcal{S}$ .

### 2 Method

#### 2.1 Rationale on Model Choice

In the previous section we identified below characteristics of the task at hand:

(1) Available Information:

(a) Known  $\mathcal{S}, \mathcal{A}$  spaces.

(b) Unknown transition probability  $P(s_{t+1} | \{s_i\}_{i=0}^t, \{a_i\}_{i=1}^t)$  and unknown reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$

(2)  $|\mathcal{A}| = 2$ , while  $\mathcal{S} \subset \mathbb{R}^6$  is continuous with  $|\mathcal{S}| = \infty$ .

(3) Outcome metric:  $\mathbb{E} \left( \sum_{s \in \mathbf{p}} R(s, \pi(s)) \middle| \mathbf{p} \right)$  the expected number of total reward in each play.

If we are willing to assume Markovian property for the process  $\{s_t\}$ , i.e.  $P(s_{t+1} | \{s_i\}_{i=0}^t, \{a_i\}_{i=1}^t) = P(s_{t+1} | s_t, a_t)$ . The available information listed in (1) put us into a Reinforcement Learning setting.

### 2.1.1 Dimension selection and Discretization

markov assumption

### 2.1.2 Exploration/Exploitation Parameters

Learning rate

$\epsilon$ -greedy

## 2.2 Performance Evaluation

## 3 Result

### 3.1 State Exploration

### 3.2 Convergence Behavior

## 4 Discussion & Possible Directions

## Reference

1. Ricci F, Rokach L, Shapira B et al. (2010) **Recommender Systems Handbook**. *Springer*.
2. Koren Y, Bell R, Volinsky C. (2009) **Matrix factorization techniques for recommender systems**. *IEEE Computer* Aug 2009, 42-49.
3. Srebro N, Jaakkola T.(2003) **Weighted low-rank approximations**. *Proceedings of the Twentieth International Conference* 720727.
4. R Salakhutdinov, A Mnih. (2008) **Probabilistic Matrix Factorization**. *Advances in Neural Information Processing Systems* Vol. 20
5. Koren, Y. (2008) **Factorization Meets the Neighborhood: a Multifaceted Collaborative Filtering Model**, *Proc. 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.