

Homework 5: Markov Decision Processes and Reinforcement Learning

Introduction

In this assignment, you will gain intuition for how MDPs and RL work. For readings, we recommend Chapters 11 and 12 of the [CS181 textbook](#) and the pre-lecture materials from April 11th and April 13th.

Please type your solutions after the corresponding problems using this L^AT_EX template, and start each problem on a new page.

Please submit the **writeup PDF to the Gradescope assignment ‘HW5’**. Remember to assign pages for each question.

Please submit your **L^AT_EX file and code files to the Gradescope assignment ‘HW5 - Supplemental’**.

You can use a **maximum of 2 late days** on this assignment. Late days will be counted based on the latest of your submissions.

Problem 1 (Markov Decision Processes, 20 pts)

In this problem, you will be working on building your conceptual understanding of MDPs.

- For the past few weeks in this course we have seen how *latent variable models* can be fit using the EM algorithm. Recall that in the M-step, we maximize the ELBO with respect to the model parameters given our best guess for the probability distribution over the latent variables. One of the terms in the ELBO is the *complete data likelihood*.

- The complete data likelihood for one latent variable model called pPCA is as follows:

$$p(Z_{1...N}, Y_{1...N} | \theta) = \prod_{n=1}^N p(Y_n, Z_n | \theta) = \prod_{n=1}^N p(Y_n | Z_n, \theta) p(Z_n)$$

Question: Use what you know from the factorization of the complete data likelihood for pPCA to either draw, or describe in words, its associated graphical model.

- The complete data likelihood for a Hidden Markov Model (HMM) is as follows:

$$\begin{aligned} p(Z_{1...N}, Y_{1...N} | \theta, \mathcal{T}) &= \prod_{n=2}^N p(Y_1, Z_1 | \theta) \prod_{n=1}^N p(Y_n, Z_n | Z_{n-1}; \theta, \mathcal{T}) \\ &= p(Y_1 | Z_1; \theta) p(Z_1) \prod_{n=2}^N p(Y_n | Z_n; \theta) p(Z_n | Z_{n-1}, \mathcal{T}) \end{aligned}$$

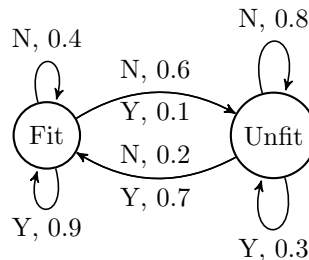
Question: How does the graphical model for the HMM differ from that for pPCA? Where do you see the Markov assumption coming into play?

- In the following example, we translate a medical study of personal fitness into an MDP. We observe the participants' performance on a number of physical activities and divide them into two categories: *fit* and *unfit*. Every week we also ask participants whether they worked out or not. We define our MDP as follows:

$$\mathcal{S} = \{\text{unfit}, \text{fit}\}$$

$$\mathcal{A} = \{\text{workout (Y)}, \text{not workout (N)}\}$$

s_n	a_n	s_{n+1}	$\mathcal{R}(s_n, a_n, s_{n+1})$
fit	Y	fit	0
fit	Y	unfit	-2
fit	N	fit	2
fit	N	unfit	-1
unfit	Y	fit	1
unfit	Y	unfit	-1
unfit	N	fit	10
unfit	N	unfit	-1



(continued on next page...)

Problem 1 (cont.)

2. Question:

- (a) What is one design choice we made in this model? What is one pro and one con of this design choice?
 - (b) Suggest one modification we could make to this model and describe what additional data (if any) we would need to collect to make this modification. How would this modification affect the rest of the model?
3. Recall that *planning* in an MDP is the process of finding a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maximizes the agent's expected reward.

Question:

- (a) In the MDP setting, why can't $\pi(s_n)$ depend on the *history* – the values of previous states s_{n-1}, s_{n-2}, \dots ?
 - (b) Is this assumption reasonable for our personal fitness MDP?
4. Recall that the *return* for a trajectory (a sequence of states and actions) of length N is given by:

$$G = \sum_{n=1}^N \gamma^n R_n$$

where R_n is the reward collected at time n .

Question: Why do we discount? Often we argue that discounting is needed if we allow infinite trajectories. Show that if $\gamma = 1$ and the trajectory is infinite then G can be undefined.

5. Recall that the value function of an MDP for state s under policy π is defined as follows:

$$V^\pi(s) = \mathbb{E}_\pi[G | Z_0 = s]$$

where the expectation above is taken over all randomly varying quantities in G .

Question:

- (a) In our personal fitness MDP what quantities that G depends on are randomly varying?
Hint: Think about sources of randomness in your reward, transition and policy.
 - (b) Suppose the agent's policy π is to exercise if in the **unfit** state and not exercise if in the **fit** state. Calculate $V^\pi(\text{fit})$ and $V^\pi(\text{unfit})$.
6. Recall that the Q-function $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{R}$ quantifies the value of a policy π starting at state s , taking action a , and *then* following policy π :

$$Q^\pi(s, a) = \mathbb{E}_\pi[G | Z_0 = s, A_0 = a]$$

Question: Assuming that the agent follows the same policy as in the previous question, calculate $Q^\pi(\text{fit}, Y)$, $Q^\pi(\text{fit}, N)$, $Q^\pi(\text{unfit}, Y)$, $Q^\pi(\text{unfit}, N)$.

Problem 2 (Reinforcement Learning, 20 pts)

In 2013, the mobile game *Flappy Bird* took the world by storm. You'll be developing a Q-learning agent to play a similar game, *Swingy Monkey* (See Figure 1a). In this game, you control a monkey that is trying to swing on vines and avoid tree trunks. You can either make him jump to a new vine, or have him swing down on the vine he's currently holding. You get points for successfully passing tree trunks without hitting them, falling off the bottom of the screen, or jumping off the top. There are some sources of randomness: the monkey's jumps are sometimes higher than others, the gaps in the trees vary vertically, the gravity varies from game to game, and the distances between the trees are different. You can play the game directly by pushing a key on the keyboard to make the monkey jump. However, your objective is to build an agent that *learns* to play on its own.

You will need to install the `pygame` module (<http://www.pygame.org/wiki/GettingStarted>).

Task: Your task is to use Q-learning to find a policy for the monkey that can navigate the trees. The implementation of the game itself is in file `SwingyMonkey.py`, along with a few files in the `res/` directory. A file called `stub.py` is the starter code for setting up your learner that interacts with the game. This is the only file you need to modify (but to speed up testing, you can comment out the animation rendering code in `SwingyMonkey.py`). You can watch a YouTube video of the staff Q-Learner playing the game at <http://youtu.be/14QjPr1uCac>. It figures out a reasonable policy in a few dozen iterations. You'll be responsible for implementing the Python function `action_callback`. The action callback will take in a dictionary that describes the current state of the game and return an action for the next time step. This will be a binary action, where 0 means to swing downward and 1 means to jump up. The dictionary you get for the state looks like this:

```
{ 'score': <current score>,
  'tree': { 'dist': <pixels to next tree trunk>,
            'top': <height of top of tree trunk gap>,
            'bot': <height of bottom of tree trunk gap> },
  'monkey': { 'vel': <current monkey y-axis speed>,
              'top': <height of top of monkey>,
              'bot': <height of bottom of monkey> }}
```

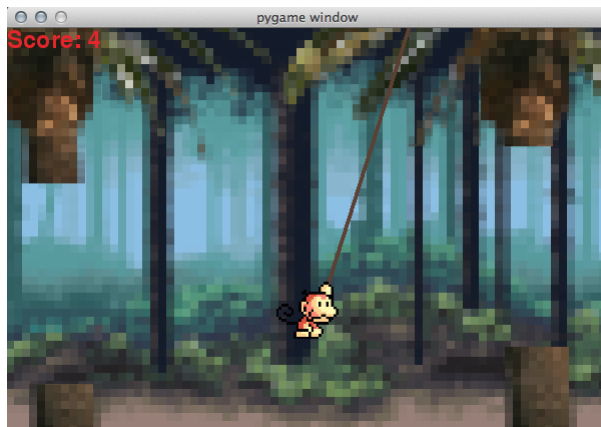
All of the units here (except score) will be in screen pixels. Figure 1b shows these graphically. Note that since the state space is very large (effectively continuous), the monkey's relative position needs to be discretized into bins. The pre-defined function `discretize_state` does this for you.

Requirements

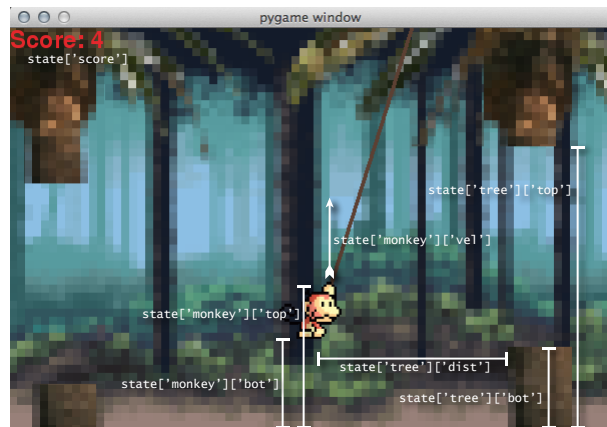
Code: First, you should implement Q-learning with an ϵ -greedy policy yourself. You can increase the performance by trying out different parameters for the learning rate α , discount rate γ , and exploration rate ϵ . *Do not use outside RL code for this assignment.* Second, you should use a method of your choice to further improve the performance. This could be inferring gravity at each epoch (the gravity varies from game to game), updating the reward function, trying decaying epsilon greedy functions, changing the features in the state space, and more. One of our staff solutions got scores over 800 before the 100th epoch, but you are only expected to reach scores over 50 before the 100th epoch. **Make sure to turn in your code!**

Evaluation: In 1-2 paragraphs, explain how your agent performed and what decisions you made and why. Make sure to provide evidence where necessary to explain your decisions. You must include in your write up at least one plot or table that details the performances of parameters tried (i.e. plots of score vs. epoch number for different parameters).

Note: Note that you can simply discretize the state and action spaces and run the Q-learning algorithm. There is no need to use complex models such as neural networks to solve this problem, but you may do so as a fun exercise.



(a) SwinglyMonkey Screenshot



(b) SwinglyMonkey State

Figure 1: (a) Screenshot of the Swingly Monkey game. (b) Interpretations of various pieces of the state dictionary.

Problem 3 (Impact Question: Assessing the energy consumption of data centers and understanding a RL approach to optimize it, 3.5 pts)

Every computation consumes energy. Over the past decades large scale data centers have been built which are significant energy consumers. Consequently, every computation is costly.

1. **Energy consumption:** How much energy do data centers consume globally per year? And what's their share in global energy consumption? (1 point)
2. **Consumption optimization with RL:** How would you formulate a RL model which enables the optimization of the load of a data center? Describe the main elements of your model (state space, action space, reward). (1.5 points)
Hint: This [RL Paper](#) provides background information for this.
3. **Alternatives:** List one alternative idea on how to reduce the climate impact of data centers. (1 point)

Name

Collaborators and Resources

Whom did you work with, and did you use any resources beyond cs181-textbook and your notes?

Calibration

Approximately how long did this homework take you to complete (in hours)?