

CS 446 / ECE 449 — Homework 6

acard6

Version 1.0

Instructions.

- Homework is due **Tuesday, April 25, at noon CST**; no late homework accepted.
- Everyone must submit individually on Gradescope under **hw6** and **hw6code**.
- The “written” submission at **hw6** **must be typed**, and submitted in any format Gradescope accepts (to be safe, submit a PDF). You may use L^AT_EX, Markdown, Google Docs, MS Word, whatever you like; but it must be typed!
- When submitting at **hw6**, Gradescope will ask you to select pages for each problem; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.
- We reserve the right to reduce the auto-graded score for **hw6code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- Coding problems come with suggested “library routines”; we include these to reduce your time fishing around APIs, but you are free to use other APIs.
- When submitting to **hw6code**, only upload the two python files **hw6.py** and **hw6_utils.py**. Don’t upload a zip file or additional files.

Version history.

2.0. Updated version.

1. Combination Lock.

Consider an MDP with n states s_1, s_2, \dots, s_n , where there are 2 actions a_1, a_2 at each state. At each state s_i , $i < n$ action a_1 takes the agent to the s_{i+1} and for state s_n , action a_1 is a self-loop which keeps the agent in s_n . At any state, action a_2 takes the agent back to s_1 , except for the last state, which again has a self-loop. See the figure for an overall picture of the setting. $R(s_i, a_j)$ (reward of action a_j at state s_i) is 0, except for (s_n, a_1) , which has a value of 1. The agent takes one action at each step.

With *uniform random policy*, the agent at each state chooses one of the available actions uniformly at random. Now considering this *combination lock* MDP, answer the questions below. You need to show your work to receive full credit for each of the sections.

- Compute the expected number of steps for the uniform random policy to go from state s_1 to state s_n .
- Compute the formula for $Q(s_i, a_j)$, $\forall i, j$ for the uniform random policy considering a discounted reward setting with a discount factor of γ .
- Prove that:

$$\forall i < n : Q(s_i, a_1) > Q(s_i, a_2).$$

- Now consider a new *greedy policy* using the $Q(s, a)$ function you computed in the previous part. Specifically, $\pi(s_i) = \operatorname{argmax}_a Q(s_i, a)$ (i.e., the agent, at each state, chooses the action with the highest value of the Q function computed in part (b)). Compute the new expected number of steps to get from state s_1 to s_k .

Hint: Use part (c); In particular, the estimate you get should have a simple form and be much smaller.

Solution.

- Lets define our recursion as $E[T|s_i] = 1 + 0.5 * E[T|s_{i+1}] + 0.5 * E[T|s_n]$, with $E[T|s_n] = 0$, since it never leaves the final state

We want to compute the expected number of steps to go from state s_1 to state s_n . To do this, we can use the recursive equation above to express $E[T|s_1]$ in terms of $E[T|s_2], E[T|s_3], \dots, E[T|s_n]$. Starting from $s_i = s_1$, we can repeatedly substitute the equation for $E[T|s_i]$ into the equation for $E[T|s_{i-1}]$, until we get an expression for $E[T|s_1]$ in terms of $E[T|s_n]$.

$$E[T|s_{n-1}] = 1 + 0.5 * E[T|s_n] + 0.5 * E[T|s_n] = 1 + E[T|s_n]$$

$$E[T|s_{n-2}] = 1 + 0.5 * E[T|s_{n-1}] + 0.5 * E[T|s_n] = 2 + 0.5 * E[T|s_n]$$

$$E[T|s_{n-3}] = 1 + 0.5 * E[T|s_{n-2}] + 0.5 * E[T|s_n] = 2.5 + 0.25 * E[T|s_n]$$

$$E[T|s_{n-4}] = 1 + 0.5 * E[T|s_{n-3}] + 0.5 * E[T|s_n] = 2.875 + 0.125 * E[T|s_n]$$

...

$$E[T|s_1] = 1 + 0.5 * E[T|s_2] + 0.5 * E[T|s_n] = \frac{(2^n - 2)}{2^{n-1}} + \left(\frac{1}{2}\right)$$

- since the $R(s_i, a_j) = 0$ except for $R(s_n, a_1) = 1$, and $Q(s, a) = (1 - \gamma) \mathbb{E} \sum_{i \geq 0} \gamma^i r_i$ this means that $Q(s_i, a_1) = (1 - \gamma) * [0 + \gamma * Q(s_{i+1}, a_1)] + \gamma * Q(s_n, a_1)$ (taking action a_1 at s_i) $= (1 - \gamma) * \gamma * Q(s_{i+1}, a_1) + \gamma * Q(s_n, a_1)$

$$Q(s_i, a_2) = (1 - \gamma) * [0 + \gamma * Q(s_1, a_2)] + \gamma * Q(s_n, a_2)$$

(taking action a_2 at s_i) $= (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma * Q(s_n, a_2)$ Note that $Q(s_n, a_1) = Q(s_n, a_2) = 1$.

We can use these equations to compute the value of $Q(s_i, a_j)$ for all i and j recursively. We start with the terminal state s_n : $Q(s_n, a_1) = Q(s_n, a_2) = 1$ Then, we can use the equations above to compute $Q(s_i, a_j)$ for $i = n - 1, n - 2, \dots, 1$:

$$Q(s_{n-1}, a_1) = (1 - \gamma) * \gamma * Q(s_n, a_1) + \gamma * Q(s_n, a_1) = 1$$

$$Q(s_{n-1}, a_2) = (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma * Q(s_n, a_2)$$

$$\begin{aligned}
Q(s_{n-2}, a_1) &= (1 - \gamma) * \gamma * Q(s_{n-1}, a_1) + \gamma * Q(s_n, a_1) = 1 \\
Q(s_{n-2}, a_2) &= (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma * Q(s_n, a_2) \\
&\dots \\
Q(s_1, a_1) &= (1 - \gamma) * \gamma * Q(s_2, a_1) + \gamma * Q(s_n, a_1) \\
Q(s_1, a_2) &= (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma * Q(s_n, a_2)
\end{aligned}$$

we can simplify the equation for $Q(s_1, a_2)$ as:

$$Q(s_1, a_2) = (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma * Q(s_n, a_2) = \gamma * Q(s_n, a_2) / (1 - (1 - \gamma) * \gamma)$$

Therefore, we have:

$$Q(s_1, a_1) = (1 - \gamma) * \gamma * Q(s_2, a_1) + \gamma * Q(s_n, a_1) Q(s_1, a_2) = \gamma * Q(s_n, a_2) / (1 - (1 - \gamma) * \gamma)$$

- (c) Using strong induction we can do it as follows.

Base case: When $i = n - 1$, we have $Q(s_{n-1}, a_1) = 1$ and $Q(s_{n-1}, a_2) < Q(s_{n-1}, a_1)$ since $Q(s_{n-1}, a_2)$ is the expected discounted reward of going from s_{n-1} to s_1 and then to s_{n-1} again, which has a maximum value of 1 (when $\gamma = 0$). Therefore, the inequality holds for $i = n - 1$.

Inductive step: Assume that the inequality holds for all states $s_{i+1}, s_{i+2}, \dots, s_{n-1}$, and consider the state s_i . We have:

$$Q(s_i, a_1) = (1 - \gamma) * \gamma * Q(s_{i+1}, a_1) + \gamma * Q(s_n, a_1) \text{ (by the equation for } Q(s_i, a_1))$$

$$Q(s_i, a_2) = (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma * Q(s_n, a_2) \text{ (by the equation for } Q(s_i, a_2))$$

By the induction hypothesis $Q(s_{i+1}, a_1) > Q(s_{i+1}, a_2)$. Also, $Q(s_n, a_1) = Q(s_n, a_2) = 1$ by the definition of the rewards

$$Q(s_i, a_1) > (1 - \gamma) * \gamma * Q(s_{i+1}, a_2) + \gamma (\text{since } Q(s_i, a_1) > Q(s_i, a_2) \text{ and } Q(s_n, a_2) = 1) \geq (1 - \gamma) * \gamma * Q(s_1, a_2) + \gamma (\text{since } Q(s_{i+1}, a_2) \geq Q(s_1, a_2) \forall i < n - 1) = Q(s_i, a_2)$$

Therefore, the inequality holds for all states $s_i, i < n$, by induction.

- (d) We note that $Q(s_1, a_1) > Q(s_1, a_2)$ since $Q(s_1, a_2)$ is the expected discounted reward of going from s_1 to s_1 repeatedly, which has a maximum value of 1 (when $\gamma = 0$). Therefore, the greedy policy will choose a_1 at s_1 .

Similarly, we have $Q(s_i, a_1) > Q(s_i, a_2)$ for all $i < n - 1$ by the previously proven inequality. Therefore, the greedy policy will choose a_1 at all states $s_i, i < n - 1$.

At s_{n-1} , we have $Q(s_{n-1}, a_2) > Q(s_{n-1}, a_1)$ since $Q(s_{n-1}, a_2)$ is the expected discounted reward of going from s_{n-1} to s_1 repeatedly, which has a maximum value of 1 (when $\gamma = 0$). Therefore, the greedy policy will choose a_2 at s_{n-1} .

Given this policy, the expected number of steps to get from state s_1 to s_{n-1} is given by:

$$E[T] = 1 + \gamma + \gamma^2 + \dots + \gamma^{n-2}$$

This is a geometric series with first term 1 and common ratio γ , so the sum can be computed as:

$$E[T] = (1 - \gamma^{n-1}) / (1 - \gamma)$$

Similarly, the expected number of steps to get from s_{n-1} to s_n is:

$$E[T'] = \frac{1}{1 - \gamma}$$

Therefore, the expected number of steps to get from s_1 to s_n using the new greedy policy is:

$$E[T_{total}] = E[T] + E[T'] = \frac{1 - \gamma^{n-1}}{1 - \gamma} + \frac{1}{1 - \gamma}$$

Simplifying this expression, we get:

$$E[T_{total}] = (1 - \gamma^{n-1}) + \frac{\gamma^{n-1}}{1 - \gamma} \text{ Therefore, the expected number of steps to get from } s_1 \text{ to } s_n \text{ using the new greedy policy is } (1 - \gamma^{n-1} + \frac{\gamma^{n-1}}{1 - \gamma} - \gamma).$$

2. Attention.

In this problem, you are given a set of binary vectors that are mapped to another set of binary vectors of the same length. You will construct convolutional and attention-based models to learn and predict this mapping. Training data and the code for training and testing your model are provided for you.

- (a) Implement a model that uses convolutional layer with the given kernel size and output channels as inputs to generating the embedding for a given sequence. After using ReLU activation function on the embedding, use a Transposed Convolution module to generate the output vector of the same length as the original input and a single channel.

Remark: Consider using `nn.Conv1d`, `nn.ConvTranspose1d`, and `torch.reshape` functions.

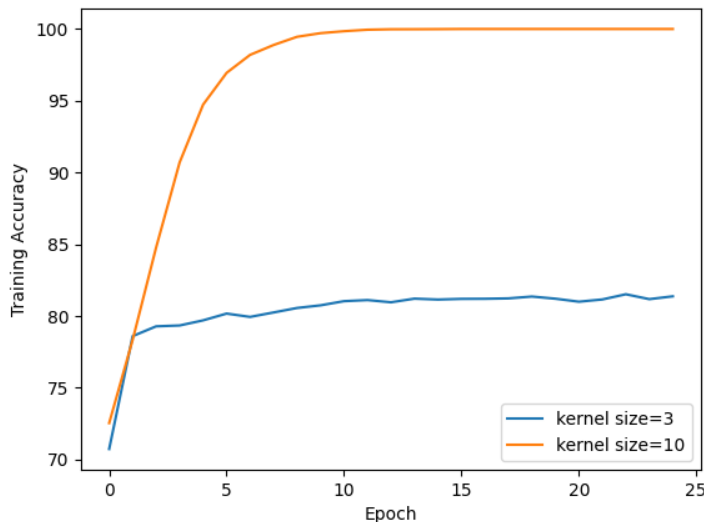
- (b) Implement an Attention-based model with a single head. The general structure is provided for you. You have to compute the queries, keys, and values using linear layers with no bias and complete the “compute_new_values” function that computes the attention matrix and the new values.

Library routines: `nn.Linear`, `F.softmax`, `torch.matmul`, and `transpose()`

- (c) Use your convolutional model with a kernel size of 10 (length of the binary sequences) and another convolutional model with a kernel size of 3. Train both for 25 epochs and plot the training accuracy curves against each other (in the same plot) for both models. Explain your observations and reason about them.
- (d) Train the attention model once with the positional encoding flag set to True (the default value) and once with this flag set to False. Plot the training accuracy curves against each other (in the same plot) for both models. Explain your observations and reason about them.
- (e) Plot the attention matrix. Can you explain in words how the sequences are mapped? Include your figure. Also, discuss how using attention and attention values are different from convolutions.

Note: Make sure not to plot the attention matrix after the model is called on the test data in the test function (you can either comment out the call to the test function or change the number of epochs to a value that the procedure does not end with a call to test function).

Solution.



(c)

(d)

(e)