

Reinforcement Learning Assignment 1

Question 1: Longest Common Subsequence

For strings $X = x_1 \dots x_n$ and $Y = y_1 \dots y_m$ of size n, m denote $lcs(X, Y)$ to be the size of the longest common subsequence. If $m = 0$ or $n = 0$ then $lcs(X, Y) = 0$. If $x_n = y_m$ then we can choose whether to take x_n and y_m to be part of the subsequence and continue with computing $lcs(x_1 \dots x_{n-1}, y_1 \dots y_{m-1})$, or we can choose to not take them as part of the sequence and calculate $lcs(X, y_1 \dots y_{m-1}), lcs(x_1 \dots x_{n-1}, Y)$. For the case where $x_n \neq y_m$ we can't take them to be part of the solution so we will calculate $lcs(X, y_1 \dots y_{m-1}), lcs(x_1 \dots x_{n-1}, Y)$.

This is concluded in the recursive formula:

$$lcs(X, Y) = \begin{cases} 0 & \text{if } n = 0 \text{ or } m = 0 \\ \max\{1 + lcs(x_1 \dots x_{n-1}, y_1 \dots y_{m-1}), lcs(X, y_1 \dots y_{m-1}), lcs(x_1 \dots x_{n-1}, Y)\} & \text{if } x_n = y_m \\ \max\{lcs(X, y_1 \dots y_{m-1}), lcs(x_1 \dots x_{n-1}, Y)\} & \text{else} \end{cases}$$

Define $M_{i,j} = lcs(x_1 \dots x_i, y_1 \dots y_j)$, the solution will be $M_{n,m}$. From the formula above we can define an iterative algorithm:

```

Input:  $X = x_1 \dots x_n, Y = y_1 \dots y_m$ 
output:  $lcs(X, Y)$ 

define  $\forall i \leq n, j \leq m \ M_{i,j} = 0$ 
For  $i = 1, \dots, n$  do:
    For  $j = 1, \dots, m$  do:
        if  $x_i = y_j$  then  $M_{i,j} = \max\{1 + M_{i-1,j-1}, M_{i-1,j}, M_{i,j-1}\}$ 
        else  $M_{i,j} = \max\{M_{i-1,j}, M_{i,j-1}\}$ 

return  $M_{n,m}$ 
    
```

Since the algorithm makes n iterations with m steps the running time is $O(nm)$.

Question 2: Moses the mouse

1) State space: $S = \{s_{i,j} | 0 \leq i \leq N, 0 \leq j \leq M\}$

Action space: $A(s_{i,j}) = \begin{cases} \{s_{i+1,j}\} & \text{if } j = M, i < N \\ \{s_{i,j+1}\} & \text{if } i = N, j < M. \\ \{s_{i+1,j}, s_{i,j+1}\} & \text{else} \end{cases}$ Meaning the action space is $\{north, east\}$ if there is no wall in the direction.

Cumulative cost function: First define the reward function

$$r_t(s_t, a_t) = \begin{cases} 1 & \text{if there is a cheese in } s_t \\ 0 & \text{else} \end{cases}$$

$$r_T(s_T) = \begin{cases} 1 & \text{if there is a cheese in } s_T \\ 0 & \text{else} \end{cases}$$

And $R(path) = \sum_{t=1}^{T-1} r_t(s_t, a_t) + r_T(s_T)$ and we can define $C(path) = -R(path)$.

2) The horizon of the problem is when Moses arrives to north-east corner $s_{N,M}$ which will happen in exactly $M + N$ steps, so $T = N + M$.

3) For $M = 2$, there must be exactly one state $s_{i,j}$ where the action *north* is taken. So each trajectory is uniquely defined by the state where Moses took the action *north*. Since there are N different states where he can take this action there are N different trajectories.

For the general case, notice that the trajectory is uniquely defined by $m - 1$ states where we choose to take the action *north*. We can think of it like choosing from $m - 1$ elements from a set of size N with repetitions. The formula for choosing k elements from n with repetitions is $\binom{n+k-1}{k}$ so for us it means that the number of different trajectories is $\binom{N+M-2}{M-1}$.

For the case where $N = M$ we get $\binom{2N-2}{N-1}$.

4) a) If both mice will go with the optimal policy, then if we assume that such policy is unique they will go through exactly the same trajectory. In this case each one will get half of the optimal reward (assuming they split the cheese). If there is more than a single optimal policy they might choose different trajectories and increase their reward.

b) For the case where the mice split and can move independently the state space will be $S \times S$ where S is the state space defined in section 1. Meaning the number of states in this case is $|S \times S| = (NM)^2$. Notice that not all the states in this set can actually be reached if we assume both mice take each step simultaneously.

Similarly, the action state is $\{\textit{north}, \textit{east}\}^2$, or more accurately $A \times A$ meaning an action consists of two legal action, one for each mice. If both *north* and *east* are legal for some state there are 4 different actions.

c) Like the section before, each mouse can choose to go *north* or *east* (if such actions are legal). So each choice of all mice corresponds to a different action, meaning the action space has size 2^K .

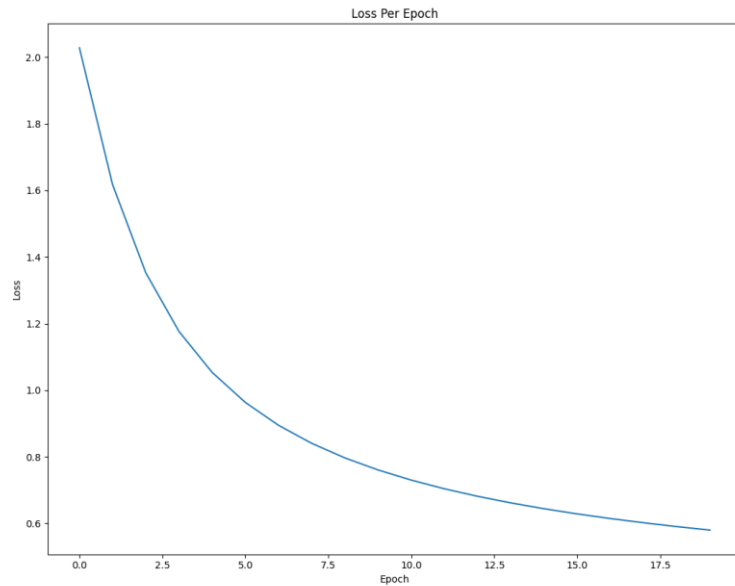
Each mouse has NM states at the room, so together K mice has state space of size $(NM)^K$.

Programing

Question 1: MNIST

1) The following plot shows accuracy per epoch for the parameters given to us:

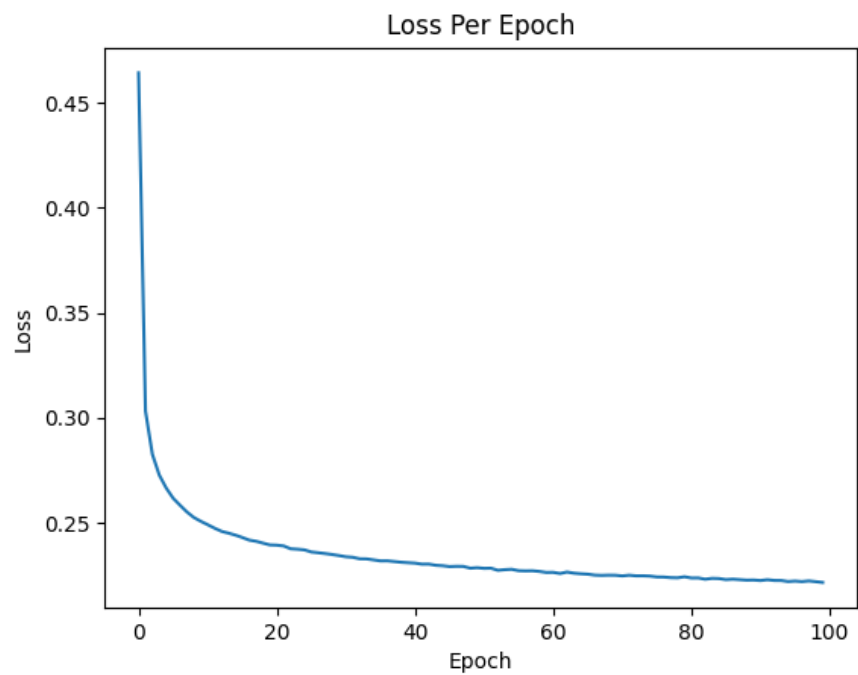
num_epoch = 20
batch_size = 100
learning_rate = 1e-3
optimizer = SGD



Accuracy of the network on the 10000 test images: 87%

2) We found the best parameters to be:

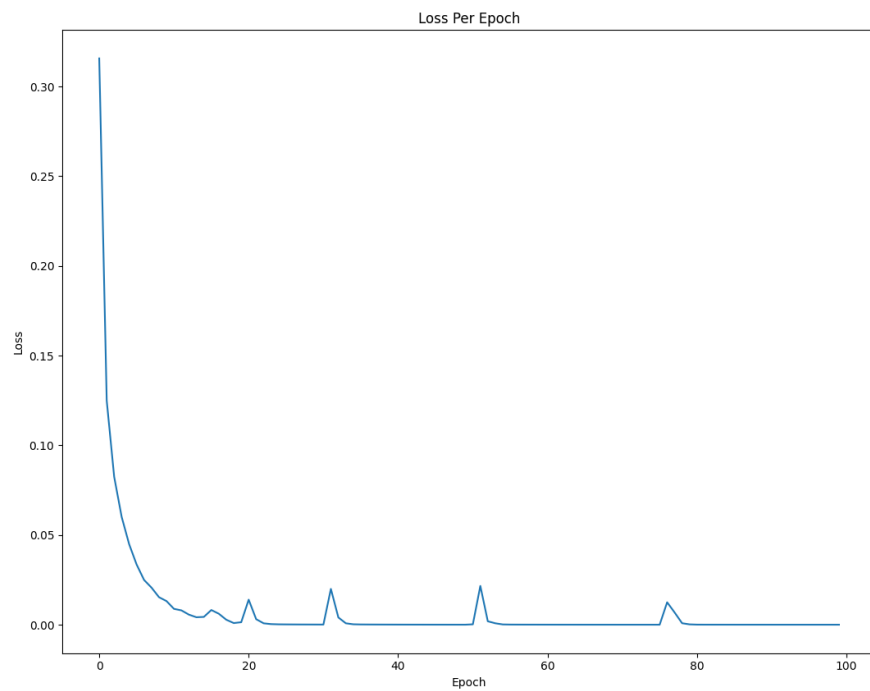
num_epoch = 100
batch_size = 32
learning_rate = 0.001
optimizer = ADAM



Accuracy of the network on the 10000 test images: 92%

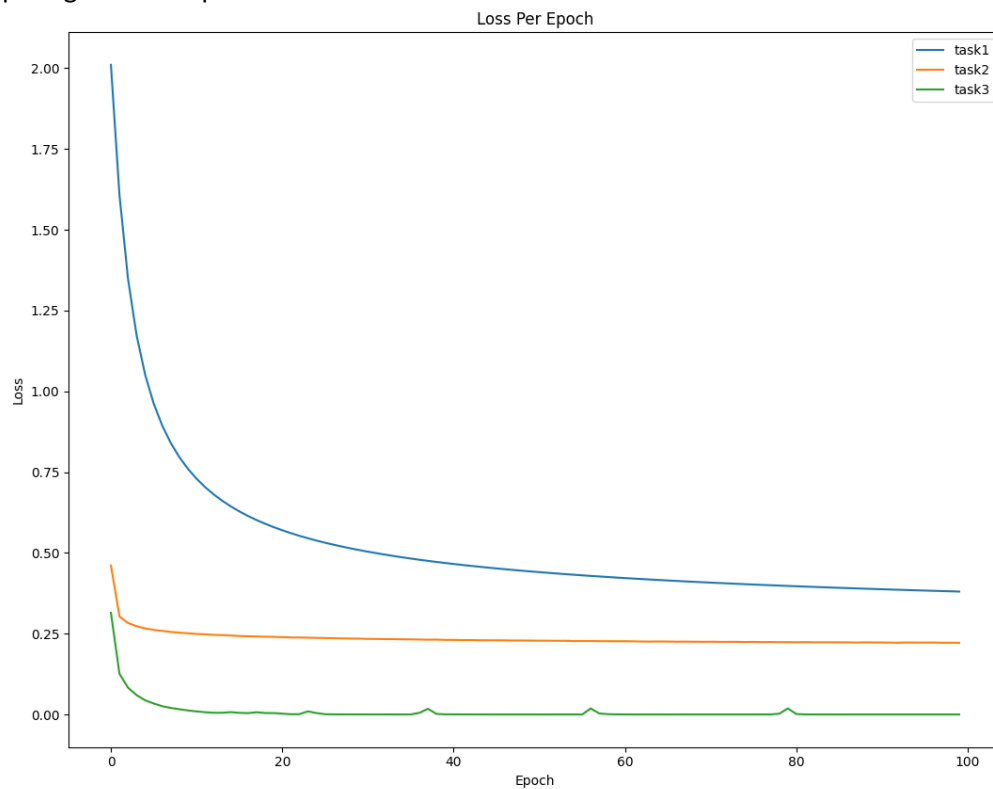
3) We found the best parameters to be:

num_epoch = 100
batch_size = 32
learning_rate = 0.001
optimizer = ADAM



Accuracy of the network on the 10000 test images: 98%

Comparing loss of all parts:



Reinforcement Learning, H.W. 1:

Q3:

$$1) P_2\{B0b\} = \frac{(b-70)}{(0.25)} \cdot \frac{(0-76)}{(0.2)} \cdot \frac{(b-7-)}{(0.325)} = 0.01625$$

(not sensory in B)

$$P_2\{OK\} = \frac{(0-7k)}{(0.2)} \cdot \frac{(k-7-)}{(0.2)} = 0.04$$

$$P_2\{B\} = 1 \leftarrow \text{CRMV was must stay in B}$$

$$P_2\{B00k\} = \frac{(b-70)}{(0.25)} \cdot \frac{(0-70)}{(0.2)} \cdot \frac{(0-7-)}{(0.2)} \cdot \frac{(k-7-)}{(0.2)} = 0.002$$

$$P_2\{B0000k\} = (0.25) \cdot (0.2) \cdot (0.4) \cdot (0.4) \cdot (0.2) \cdot (0.2) = 0.0008$$

! (11R)

✓

1

(2) a)

↓

• State-space:

$$\left\{ \begin{array}{l} S_0 := \{ (b') \} \\ S_{1 \leq t \leq k-2} := \{ (b'), (k'), (o') \} \\ S_{k-1} := \{ (i') \} \end{array} \right.$$
$$\Rightarrow S = \bigcup_{t=0}^{k-1} S_t$$

• Action-space:

$$\left\{ \begin{array}{l} A_{0 \leq t \leq k-3} := \{ (b'), (k'), (o') \} \\ A_{k-2} := \{ (i'-') \} \end{array} \right.$$
$$\Rightarrow A = \bigcup_{t=0}^{k-2} A_t$$

• Deterministic transition-function:

$$\left\{ \begin{array}{l} f_{0 \leq t \leq k-3} (s_t, a_t) = a_t \\ f_{k-2} (s_{k-2}, (i'-')) = (i'-') \end{array} \right.$$

↓

(2)

• multiplicative cost-function:

$$L_t(s_t, a_t) = 1 - P_t(s_t - \gamma a_t)$$

!important! of

\Rightarrow overall-cost:

$$\prod_{t=0}^{K-1} L_t(s_t, a_t)$$

$V(a, \hat{s})$

(b)

לחשב $L_t(s)$ נדרש

עד $t=K-1$ $K-1$ מסתכל על t מסתכל s - נדרש $from$ pe

DP of L_t

- Init: $L_{K-1}(s) = 1$; $\forall s \in S_{K-1}$

- backward-recursion:

for $t = 0, \dots, K-2$:

$$L_t(s) = \min_{a \in A_t} \{ L_t(s, a) \cdot L_{t+1}(f_t(s, a)) \}$$

$$\Pi_t(s) = \arg \min_{a \in A_t} \{ L_t(s, a) \cdot L_{t+1}(f_t(s, a)) \}$$

\Downarrow 2

↓

• כלומר, נרצה למצוא את \mathcal{H}_t ו- \mathcal{H}_t (ב')
 נראה כי \mathcal{H}_t הוא מרחב וקטורי.
 • \mathcal{H}_t הוא מרחב וקטורי.

- נשים לב, כי \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.
 • \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.
 • \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.

! (ב')

(ג)
 • נשים לב, כי \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.

$$\log \left(\prod_{t=0}^{k-1} \mathcal{H}_t(s_t, a_t) \right) = \sum_{t=0}^{k-1} \log (\mathcal{H}_t(s_t, a_t))$$

 • נשים לב, כי \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.

! (ג')

(ד)
 - נשים לב, כי \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.
 • \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.
 • \mathcal{H}_t הוא מרחב וקטורי, ו- \mathcal{H}_t הוא מרחב וקטורי.

! (ד')

Bk

14

Intro RL – EX1

Theoretical section:

Question 4:

1. The following is the results of calculating $d_i(v)$ for every $i \in \{0, \dots, n = 5\}, v \in \{A, B, C, D, E\}$ given initial state $s = A$:

	d_0	d_1	d_2	d_3	d_4	d_5	$\max(\frac{d_n(v)-d_k(v)}{n-k})$
A	0	∞	∞	6	10	7	7/5
B	∞	-1	∞	2	-1	5	6
C	∞	∞	2	6	3	2	0
D	∞	∞	1	-2	4	1	3/2
E	∞	∞	1	∞	4	1	0

2. According to Karp's theorem the cost of the minimum mean cost cycle is:

$$u^* = \min_{v \in V} \max_{0 \leq k \leq n-1} \frac{d_n(v) - d_k(v)}{n - k} = 0$$

This is obtained for either $(v = C, k = 2)$ or $(v = E, k = 2)$.

3. As was proven in lecture 2, the optimal average cost of a DDP equals exactly to the cycle with the lowest average cost, hence it is 0. The cycle is $B \rightarrow E \rightarrow D \rightarrow B$.

Programming section:

Question 2:

4. The average number of required random initializations is: 13.491.
Below is the histogram:

