

CS4122: Reinforcement Learning

Programming Assignment 1

SE22UARI165 | SE22UARI123 | SE22UARI054

CONTRIBUTIONS of each team member:

SRI SPOORTHY VATTEM (SE22UARI165):

- Question 1
- Question 5 Policy Gradient (policy_gradient.py)
- Question 6 (tested for policy gradient)

JOTSNA SREE PAPPU (SE22UARI123):

- Question 4 UCB(lin_ucb.py)
- Question 6 (tested for UCB)
- Question 7 (b),(c)

GUDURI SRIKRUTH VARMA (SE22UARI054):

- Question 3 ϵ -greedy(lin_greedy.py)
- Question 2
- Question 6 (tested for ϵ -greedy)
- Question 7 (a)

Section 7 The Deliverables

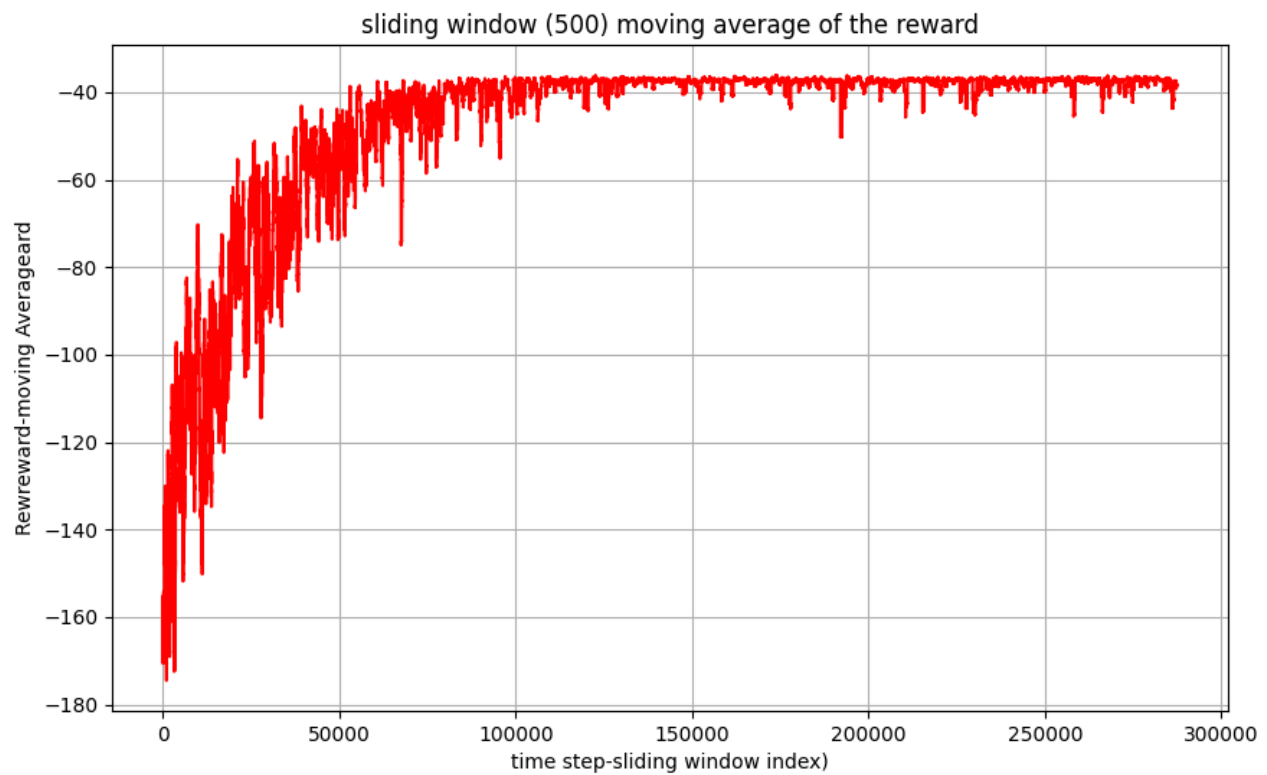
1. Contextual bandits deal with stochastic setups, in ServerAllocationENV the sources of noise include:
 - Number of jobs arriving at data center for processing (n_t) : no. of jobs change from one time slot to another and is a random integer between 1 & 8. Size of context is affected by this as no. of jobs keep changing per each time slot.
 - Cost incurred in every time slot as it depends on several factors like (wait time of all the jobs in the batch and energy cost of running the allocated servers etc..) even if the no. of jobs in different time slots are same the jobs might come with different job features like priority, latency, job type.
 - The data center knows the estimated processing time of each job but doesn't know the true processing time which introduces some randomness or unpredictable fluctuations in the environment. Latency of the job is a random variable for the data center.
 - The type of job, as ML related jobs require huge computational resources, some jobs require high speed processing and some are I/O intensive jobs and the network usage varies.
 - Overall impact on the reward, due to the above mentioned sources of randomness, the reward which is the negative of the cost is noisy which is required for this setup and helps the agent learn in different conditions.
2. (a) How to deal with the fact that the context size is varying?
 - For this problem every time a new batch of jobs arrives the number of jobs can change . so, the context also changes in the size of, because the more jobs there are the more data points we get. In most machine learning models we expect the input to be the same size. To fix this we can take all the information from jobs and average it into a single , fixed-size vector. For example: instead of looking at each Job's priority like type, network and processing time individually, we can calculate the average of these values for all jobs in that time slot. This gives us a consistent set of numbers regardless of how many jobs there are. This averaged data can then be fed into the RL agent for training.
 - There are other ways to deal with this too like using :
Padding: extend shorter sequences with a placeholder.
Truncation: shorten longer sequences to a fixed length.
But for this problem, averaging is an easy and effective solution from my point of view.
- (b) Can we reduce the number of features?

YES, we can reduce the number of features to make the problem much easier for the agent to learn from.

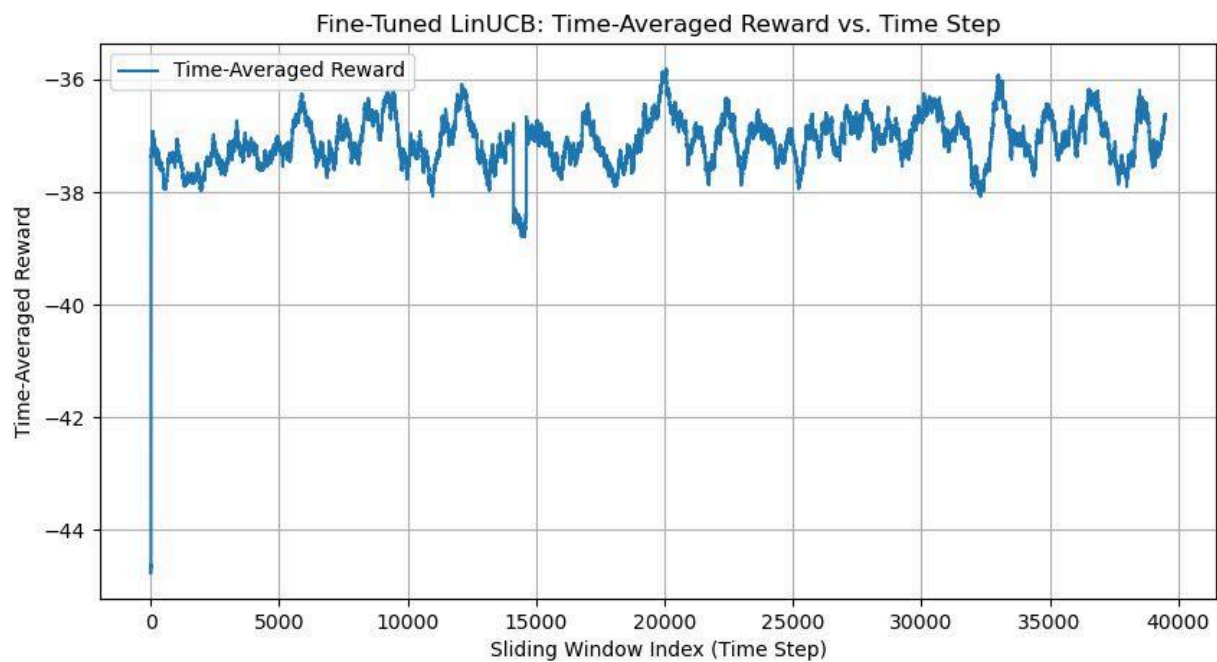
Each job has four pieces of information: Priority, job type, Network usage, Estimated processing

 - cross check if some features are related:
If two features like network usage and processing time always seem to go up and down together, we could combine them or remove one of them.
 - use techniques to simplify the data:
We can use methods like Principal Component Analysis to combine features in a way that reduces the total number of them but keeps the most important information.
We can even use Mutual Information for feature selection.
 - create new and combined features:
we could combine job priority and processing time into one value if they both affect how long the jobs wait.

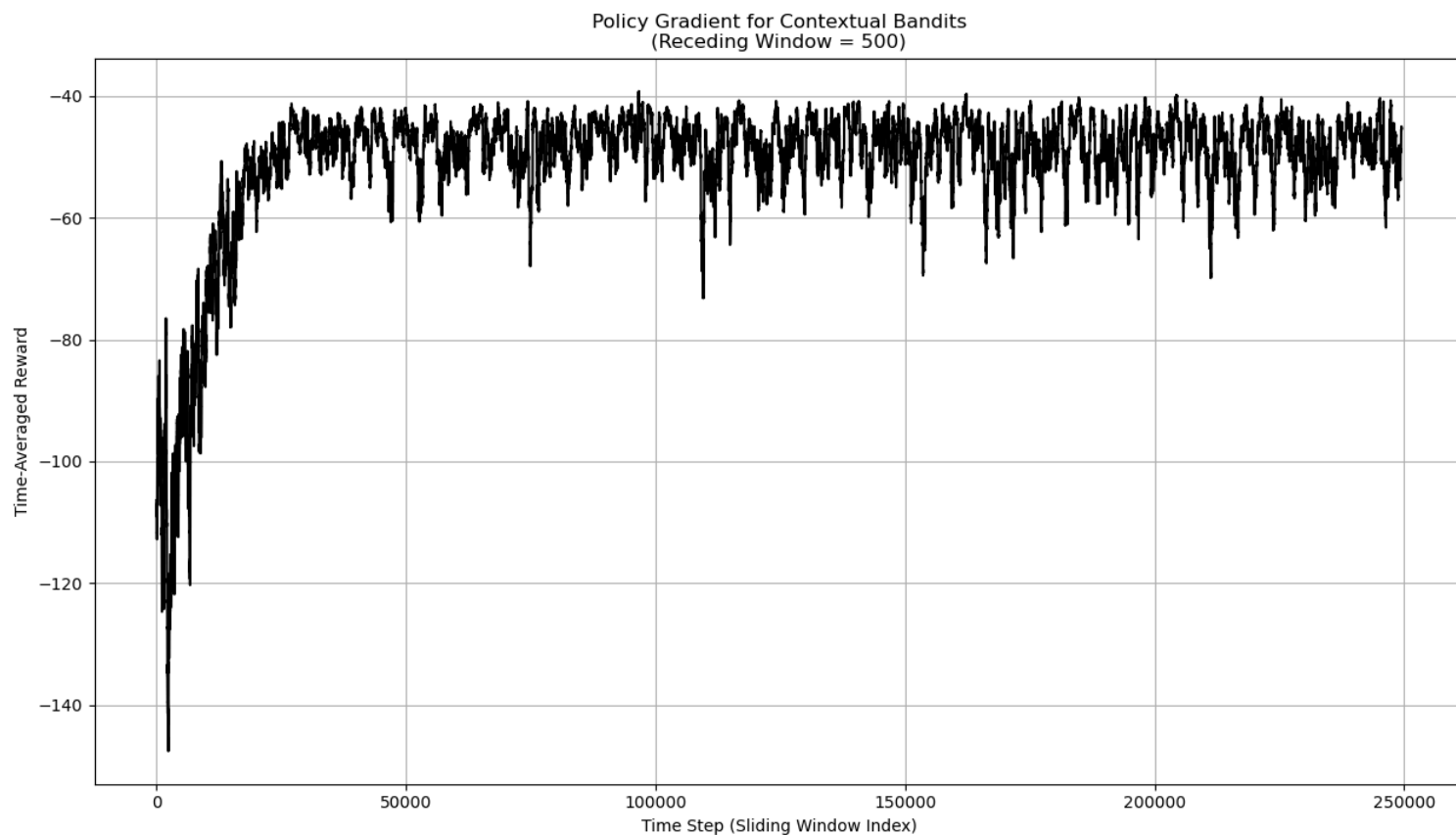
3. E-greedy graph-



4. UCB graph- Jotsna

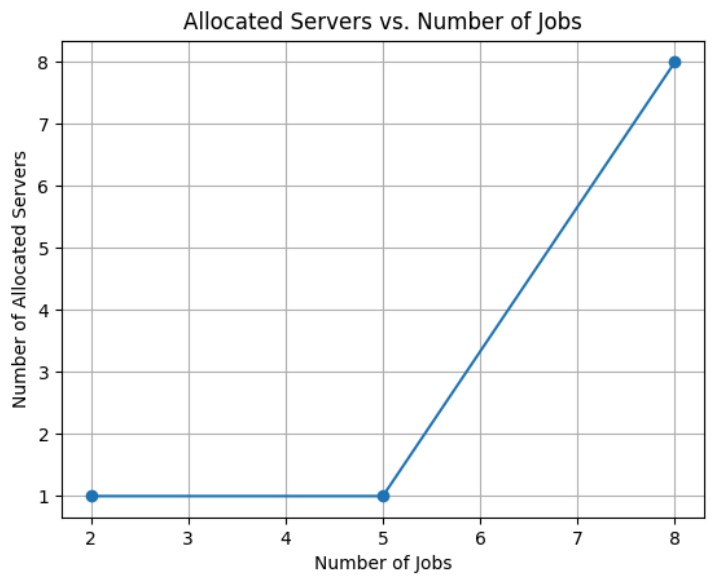
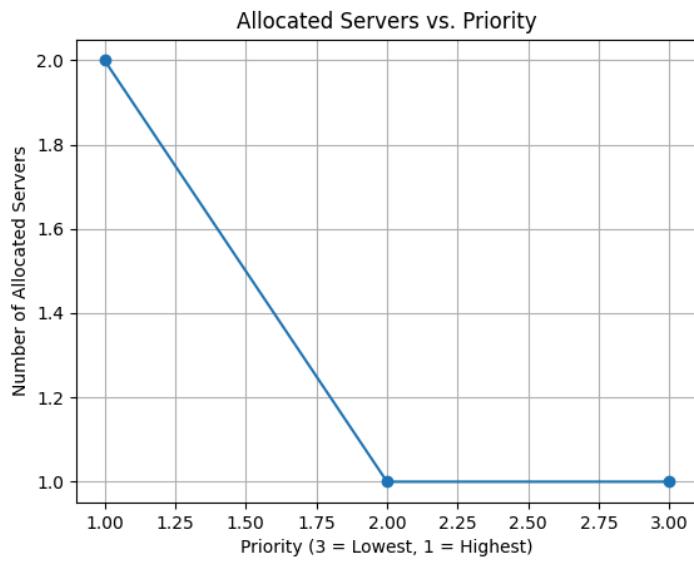


5. Policy Gradient plot:

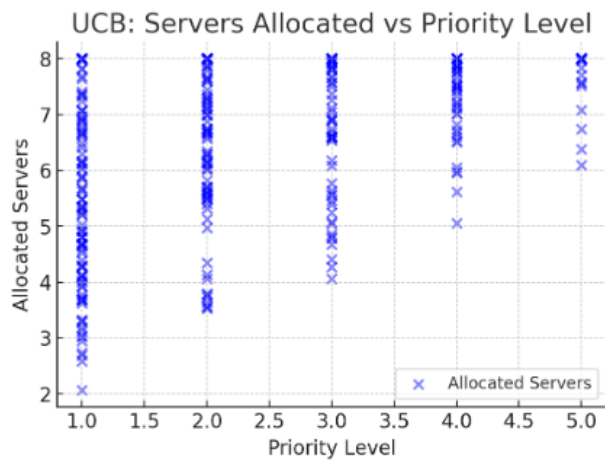


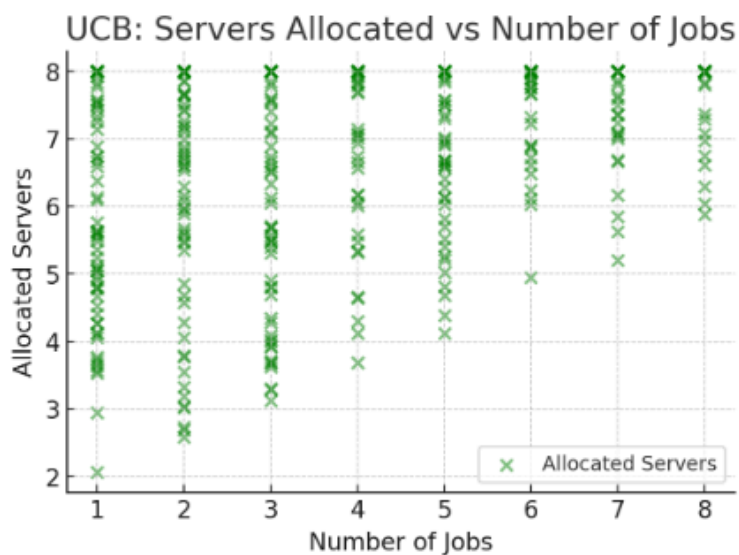
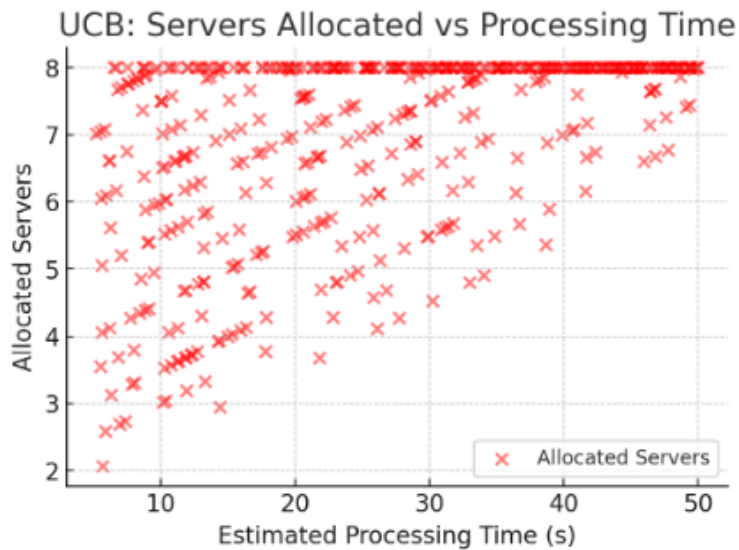
6. Testing the policies:

- ϵ -greedy:



- UCB: [Jotsna]





UCB Logical Evaluation Results

① Servers Allocated vs Priority Level

- The plot shows a **positive correlation**: as **priority increases**, the allocated servers **increase**.
- This indicates that UCB is correctly prioritizing **high-priority jobs**.

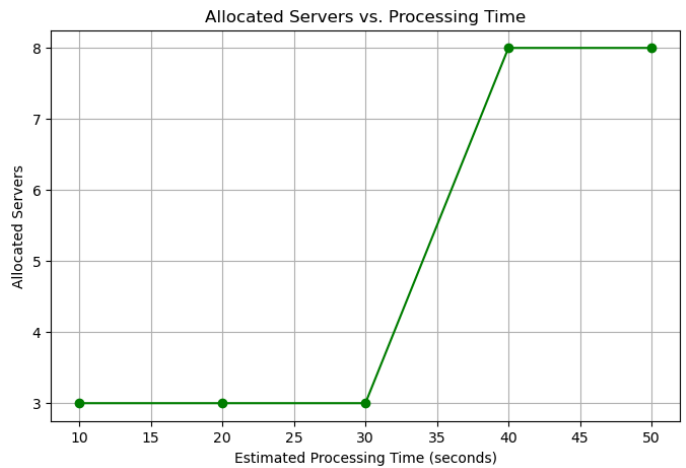
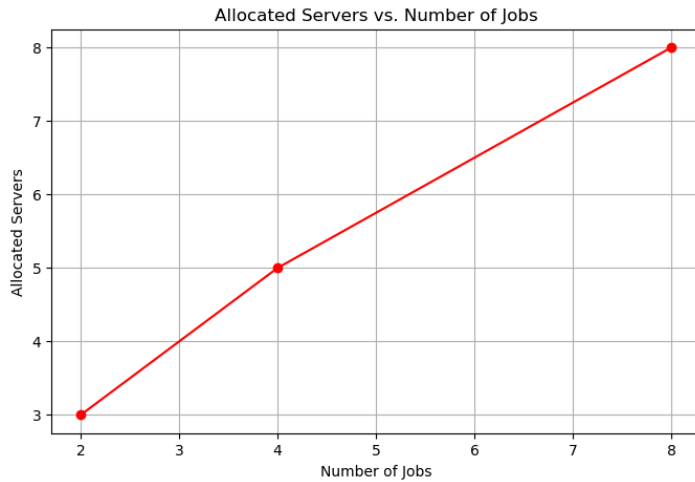
② Servers Allocated vs Estimated Processing Time

- The allocation **increases with processing time**.
- This suggests that the UCB agent **allocates more servers to longer jobs**, which is **expected behavior**.

③ Servers Allocated vs Number of Jobs

- The plot shows a **positive trend**: as **number of jobs increases**, more servers are allocated.
- This confirms that the UCB agent **adjusts allocations dynamically**.

- Policy gradient:



7. (a) Gaussian distribution
import numpy as np

```
mean = np.dot(theta1.T,x)
std_dev = np.exp(np.dot(theta2.T,x))
action = np.random.normal(mean, std_dev)
```

(b)

• update θ_1 and θ_2 using Gradient Ascent.

the parameter of the gaussian policy: $\mu, \sigma \rightarrow \theta$

$$\pi(a|x, \theta) = \frac{1}{\exp(\theta_2 x) \sqrt{2\pi}} \exp\left(-\frac{(a - \theta_1 x)^2}{2(\exp(\theta_2 x))^2}\right)$$

Step 2:

the likelihood function taking the log of the policy

$$\begin{aligned} \log \pi(a|x, \theta) &= -\log(\exp(\theta_2 x) \sqrt{2\pi}) - \frac{(a - \theta_1 x)^2}{2(\exp(\theta_2 x))^2} \\ &= -\theta_2 x - \frac{1}{2} \log(2\pi) - \frac{(a - \theta_1 x)^2}{2e^{2\theta_2 x}} \end{aligned}$$

Step 3:

gradient w.r.t θ_1 ,

$$\frac{\partial}{\partial \theta_1} \log \pi(a|x, \theta) = \frac{(a - \theta_1 x)x}{e^{2\theta_2 x}}$$

using policy gradient estimation:

$$\hat{V}_{\theta_1} J(\theta) = \frac{1}{N} \sum_{t=1}^N r_t \frac{(a_t - \theta_1 x_t)x_t}{e^{2\theta_2 x_t}}$$

Step 4: computing gradient w.r.t θ_2 .

$$\frac{\partial}{\partial \theta_2} \log \pi(a|x, \theta) = -x + \frac{(a - \theta_1 x)^2 x}{e^{2\theta_2 x}}$$

$$\hat{V}_{\theta_2} J(\theta) = \frac{1}{N} \sum_{t=1}^N r_t \left(-x_t + \frac{(a_t - \theta_1 x_t)^2 x_t}{e^{2\theta_2 x_t}} \right)$$

Step 5: update rules (Gradient Ascent)

$$\theta_1 \leftarrow \theta_1 + \eta \frac{1}{N} \sum_{t=1}^N r_t \frac{(a_t - \theta_1 x_t)x_t}{e^{2\theta_2 x_t}}$$

$$\theta_2 \leftarrow \theta_2 + \eta \frac{1}{N} \sum_{t=1}^N r_t \left(-x_t + \frac{(a_t - \theta_1 x_t)^2 x_t}{e^{2\theta_2 x_t}} \right)$$

where η is learning rate

(c) Policy Gradient Pseudocode

- For episode = 1, ..., E:

- Create an empty dataset `episode_training_data`
- For $t=1,\dots,N$:
 - Sample context x_t from context distribution F_e .
 - Sample action $a_t \sim N(\theta_1 x_t, \exp(2\theta_2 x_t))$
 - Execute a_t and observe reward r_t
 - Append (x_t, a_t, r_t) onto `episode_training_data`

Use `episode_training_data` to update θ using stochastic gradient ascent:

- Compute policy gradients:

$$\begin{aligned}
 \bullet \quad \nabla_{\theta_1} &= \frac{1}{N} \sum_{t=1}^N r_t \frac{(a_t - \theta_1 x_t) x_t}{e^{2\theta_2 x_t}} \\
 \bullet \quad \nabla_{\theta_2} &= \frac{1}{N} \sum_{t=1}^N r_t \left(-x_t + \frac{(a_t - \theta_1 x_t)^2 x_t}{e^{2\theta_2 x_t}} \right)
 \end{aligned}$$

- Update parameters using gradient ascent:

- $\theta_1 \leftarrow \theta_1 + \eta \nabla_{\theta_1}$
- $\theta_2 \leftarrow \theta_2 + \eta \nabla_{\theta_2}$