

homework-5-ishagore-1

April 14, 2024

1 NEU502B Homework 5

Due April 15, 2024

Submission instructions: First, rename your homework notebook to include your name (e.g. `homework-5-nastase.ipynb`); keep your homework notebook in the `homework` directory of your clone of the class repository. Prior to submitting, restart the kernel and run all cells (see *Kernel > Restart Kernel and Run All Cells...*) to make sure your code runs and the figures render properly. Only include cells with necessary code or answers; don't include extra cells used for troubleshooting. To submit, `git add`, `git commit`, and `git push` your homework to your fork of the class repository, then make a pull request on GitHub to sync your homework into the class repository.

In the first homework assignment, we explored how a system can extract latent structure in sensory stimuli (e.g. natural scenes) using unsupervised learning algorithms like Hebbian learning. Our model was shown a set of images with no final goal specified, nor any expectations with which to compare its performance throughout learning. Now, we're interested in how a system can learn to reach a goal through interactions with its environment, by maximizing rewards or minimizing penalties.

Reinforcement learning (RL) models solve problems by maximizing some operationalization of reward. These models use goal-directed learning to solve closed-loop problems: present actions influence the environment, thus changing the circumstances of future actions toward the same goal. In RL, we hope to discover the actions that increase chances of rewards within specific states in the environment.

An agent must be able to sense the state of the environment either fully or partially, and its actions must be able to change this state. Consider the following example from Sutton and Barto (1992):
> “Phil prepares his breakfast. Closely examined, even this apparently mundane activity reveals a complex web of conditional behavior and interlocking goal-subgoal relationships: walking to the cupboard, opening it, selecting a cereal box, then reaching for, grasping, and retrieving the box. Other complex, tuned, interactive sequences of behavior are required to obtain a bowl, spoon, and milk jug. Each step involves a series of eye movements to obtain information and to guide reaching and locomotion. Rapid judgments are continually made about how to carry the objects or whether it is better to ferry some of them to the dining table before obtaining others. Each step is guided by goals, such as grasping a spoon or getting to the refrigerator, and is in service of other goals, such as having the spoon to eat with once the cereal is prepared and ultimately obtaining nourishment.”

At each point in time, there is a state-action pair. Some of them fall under sub-goals, while others could ultimately be a state where there is a high chance of reward, fulfilling the goal of feeding. To be able to model this process, we have to break it down into its interacting components: - The

agent has a policy, the map between perceived states and the actions taken. We can think of it as a set of stimulus-response rules or associations that determine behavior given a state and a goal within the environment. It can be implemented through the probabilities of taking specific actions given a state. - This set of rules should serve to maximize the reward signal in the short and/or long term. - Environmental states are evaluated through a value function, which provides a measure of the expected rewards that can be obtained moving forward from a specific state. Grabbing a bowl might not feed you immediately, yet it has high value as it will lead you to a state in which you can feed yourself some cereal without spilling milk all over the table. Would grabbing a shallow dish instead of a bowl have the same value? Actions are taken based on these value judgements. - The agent could have the ability for foresight and planning if it has a model of the environment. This means it can have a model of how the environment reacts to its behavior, from which to base its strategies and adjustments.

At each decision, the agent has a choice to either exploit the actions it has already tested to be effective, or it can explore the action-state space to find new routes to optimal rewards. Exploration is risky, yet under some circumstances it will pay off in the long run. Finding the balance between the two would be the optimal solution in uncertain environments. Different methods can be employed to deal with this duality: - On-policy methods improve the policy that is used to make decisions. This policy is generally soft (probabilistic), as $P(s \in S, a \in A | s) > 0$, where S is the possible states and $A|s$ is the possible actions given a state. The probability is gradually shifted to a deterministic optimal policy with each update. For example, ϵ - *greedy* policies choose an action that has maximal expected value most of the time (with probability $1 - \epsilon$ a small number ϵ). However, with probability ϵ the agent will choose an action at random. The agent will try to learn values based on subsequent optimal behavior, yet it has to behave non-optimally (choosing random actions) in order to explore and find the optimal actions. This means the agent has to learn about the optimal policy while behaving according to an exploratory policy. On-policy can be thought of as a compromise, where values are learned for a near optimal policy that still explores. - Another approach is to use two policies, a target policy and a behavior policy. The first one is modified based on the information that is collected through the behaviors generated by the second. This approach is termed off-policy, as learning occurs based on behaviors generated off the policy being developed. The benefit here is that the target policy can be deterministic (i.e. greedy), while the behavior policy can continue to explore without limits.

```
[3]: import numpy as np
import matplotlib.pyplot as plt
from scipy.optimize import minimize
```

1.0.1 Problem 1: Fitting RL models to data

First, familiarize yourself with the two-step RL task (Daw et al., 2011). Visit this website to play through an example of the two-step RL task: <https://nivlab.github.io/jspsych-demos/tasks/two-step/experiment.html>. If you're interested, the Python code for the task can be found at: <https://github.com/nivlab/jspsych-demos/tree/main/tasks/two-step>.

The data from the two-step task are structured as follows: - **choice1**: your choices at the first level (1 or 2) - **choice2**: your choices at the second level (1 or 2) - **state**: which second level game you were offered on this trial - choice1 = 1 at the first level (S1) leads to S2 in approximately 70% of the trials - choice1 = 2 at the first level (S1) leads to S3 in approximately 70% of the trials - **money**: did you get a reward on each trial or not (0 or 1)

Note that missed trials will have a 0 in the choice; trials can be missed either at the first or second level. When you write your code (later on), make sure to deal separately with missed trials as this is a common source of discrepancies while fitting the models. Below is a schematic representation of the task structure:

	A_1	A_2
S_1	$Q_{S1,A1}$	$Q_{S1,A2}$
S_2	$Q_{S2,A1}$	$Q_{S2,A2}$
S_3	$Q_{S3,A1}$	$Q_{S3,A2}$

The schematic does not map to the colors used in the actual task. S_1 refers to the state at the top (first) level, where you will be shown two distinct rocket ships. You will have to choose one of the two (represented by action | state in the schematic). One of the rockets, let's say $A_1|S_1$ will have 70% chance of transferring you to S_2 (one of the possible states at the bottom level), and a 30% chance of getting you to S_3 . This is represented by the thickness of the arrows. For $A_2|S_1$, the chances are inverted. At the bottom (second) level, you can be at either of two distinct states (S_2 or S_3). You will need to choose between two aliens at each state, with gradually drifting chances of getting a reward once a decision is made. For example, $A_1|S_2$ might start with higher chances than $A_2|S_2$. These probabilities will change gradually with time and at some point, the chances might be reversed. There is no implicit relationship between what happens in S_2 and S_3 . You will have to learn, with each experience, which choices lead you to better rewards.

Let's load in the data for one of our subjects. Make sure you understand what each variable contains.

```
[4]: data = np.load('sub-0.npz')

c1 = data['choice1']
c2 = data['choice2']
s = data['state']
m = data['money']
```

At the bottom level, learning can be modeled with Q -learning or Rescorla-Wagner learning, as there's no future state. Note that these learning rules are identical if you treat each option as an action (in Q -learning) or as a state (the state of the chosen stimulus, in Rescorla-Wagner).

Q -learning: $Q^{new}(a|s) \leftarrow Q(a|s) + \eta * (R_t - Q(a|s))$

Rescorla-Wagner learning: $V_{t+1} \leftarrow V_T + \eta * (R_T - V_T)$

Question: Describe in words the variables in each equation how each the learning rule “works”:

[your answer here] For the Q -learning: $Q^{new}(a|s)$: This is the updated value estimate for taking action a in state s . $Q(a|s)$: This is the current value estimate before the update. η : This is the learning rate, a parameter that determines how much new information affects the current value estimate. A higher learning rate means that new information is weighted more heavily. R_t : This is the actual reward received at time t . $(R_t - Q(a|s))$: This is the reward prediction error, the difference between the expected reward for action a in state s and the reward actually received. The Q -learning rule

updates the value estimate for the chosen action by moving it in the direction of the reward prediction error, scaled by the learning

For Rescorla-Wagner learning: V_{t+1} : This is the updated value estimate for the stimulus at time $t + 1$. V_T : This is the current value estimate for the stimulus at time T . η : This is the learning rate. R_T : This is the actual reward received at time T . $(R_T - V_T)$: This is the reward prediction error for the stimulus

Both learning rules rely on predictions based on the reward prediction error. By updating the value estimates or Q values after each trial, the learner refines their predictions to more accurately anticipate future rewards based on their actions.

At the top level, learning can be modeled in several different ways. We'll consider two: (1) **model-free learning** and (2) **model-based learning**.

For **model-free learning**, we'll start with the temporal difference (TD) learning rule.

$$\text{TD}(0): V_T \leftarrow V_T + \eta * (R_T + \gamma * V_{T+1} - V_T)$$

Here, $R_T = 0$ because the first state doesn't yield rewards and γ is the temporal discount parameter of future rewards—this allows us to adjust the first-stage actions by taking into account the result of the second-stage action.

One way to make learning more efficient is to use TD(λ) instead of TD(0) learning. In this case, we add an additional memory variable associated with each state to serve as an “eligibility trace”. You can think of it as a “memory” that a particular state has been visited, which decays (e.g. exponentially) over time. Every time a state is visited, its eligibility trace becomes 1; at every subsequent time point, the eligibility trace is multiplied by a factor $0 < \lambda \leq 1$. At the end of a trial or episode, all eligibility traces become 0.

All states are updated according to *learning rate · prediction error · eligibility trace*. This will automatically update all the states visited in this episode (i.e. all the states “eligible” for updating), doing so for the most recently visited states to a greater extent. Write the updated equation.

$$\text{TD}(\lambda): V_{T+1} = V_T + \eta * E(\lambda) * (R_T + \gamma V_{T+1} - V_T)$$

Question: Again, describe the variables in these equations and how the learning rules “work”:

[your answer here] For TD(0): V_T : The estimated value of the current state at time T . R_T : The reward received after transitioning from the current state (which is zero for the first state). V_{T+1} : The estimated value of the next state at time $T + 1$. γ : The discount factor for future rewards, which determines the importance of future rewards. A lower value places more emphasis on immediate rewards, while a higher value favors long-term rewards. η : The learning rate. The update rule for TD(0) reflects the idea that the value of the current state should be closer to the reward received plus the discounted value of the next state

TD(λ) is similar to TD(0) with an addition of eligibility traces. $E(\lambda)$: This is the eligibility trace for a given state, which decays over time by a factor of λ each step. When a state is visited, the trace is set to 1, and it decays at each time step unless the state is visited again. λ : This is the trace-decay parameter, which determines the rate at which the memory of visiting a state decays. It is between 0 and 1, with higher values indicating a slower decay. The learning rule for TD(λ) means that each state's

value is updated not just based on the latest experience, but also taking into account how recently and frequently the state has been visited. *tate*.

For **model-based learning**, let's begin by assuming that transition model (i.e. the probabilities of going from $S1$ to $S2$ or $S3$ given choice1) is known from the start—while the reward model is not known.

Question: How can you use the transition probabilities and the learned values at the second-stage states to plan and make choices at the first stage? How would you implement this model?

[your answer here] The agent can learn transition probabilities and learned values at second stage to make choices in the first stage using the following strategy: through repeated interactions with the environment, the agent must use observed rewards to learn the values of actions at the second stage using the Q -learning or Rescorla-Wagner rule. Then, use the transition model to calculate the expected value of first-stage actions. The expected value of a first-stage action is the sum of the values of the second-stage states weighted by the transition probabilities. Choose the action with the highest expected value calculated in the previous step. As the agent continues to interact with the environment, it will gain more information about the rewards associated with different actions at the second stage, which can be used to refine the estimated values of $S2$ and $S3$, leading to better decision-making at $S1$. This process will continue iteratively, with the agent continuously refining its value estimates and choices based on new experiences, eventually converging on a strategy that maximizes rewards.

Question: How many parameters do each the four models have?

[your answer here] The Q model has one parameter which is the learning rate along with Q value for each action-state pair

The Rescorla-Wagner model just like the Q model has one parameter that is the learning rate along with V value for each stimulus.

In the $TD(0)$ model we have 2 parameters, the learning rate and gamma, the discount factor along with V value for each state.

In the $TD(\lambda)$ model we have 3 parameters, the learning rate, discount factor and trace decay parameter along with V value for each state

Now we'll implement and fit the following models. Implement $TD(\lambda)$ using the Q -learning and State-Action-Reward-State-Action (SARSA) algorithms. Some pseudocode is provided to get you started. These algorithms use state-action value predictions (Q values) to choose actions. In state S , the algorithm chooses an action according to softmax Q values.

Here, β is an inverse-temperature parameter that we'll optimize. If you're using constrained optimization, fix β to be in the range $[0, 100]$.

Update the eligibility traces. Recall that the eligibility traces are values corresponding to each state and action pair, and are set to zero at the beginning of the trial. Upon taking action a to leave state S for state S^{new} and receiving reward r , the eligibility traces $e(a | S)$ are updated for each (S, a) pair:

All $Q(a|S)$ are updated according to:

With prediction error $\delta(t)$ being:

The parameter η is a step-size or learning-rate parameter in the range (0,1].

Reset the eligibility traces to 0 at the end of each round.

```
[6]: def rl_nll(eta, beta, lambd, state, choice1, choice2, money, model_based=False):  
    """  
  
    """  
    eta, beta, lambd = params  
    n_states = int((np.max(state)) + 1)  
    n_actions = 2  
    n_trials = len(state)  
  
    # Initialize an array to store Q-values that is size n_states x n_actions  
    Q = np.zeros((n_states, n_actions))  
  
    # Initialize log-likelihood  
    LL = 0  
  
    for t in range(n_trials):  
  
        # Create an n_states x n_actions matrix to store your eligibility_  
        ↪traces for the current trial  
        E = np.zeros((n_states, n_actions))  
  
        # Get your current state for the top level (S1)  
        S = 0  
  
        # Stop if trial was missed. Missed trials will have a value of -1.  
        if choice1[t] == -1:  
            continue  
  
        # First level choice likelihood: compute likelihood of choice at the_  
        ↪first state S1.  
        # Your likelihood should be a softmax function.  
        q_values = Q[S]  
        exp_q = np.exp(beta * q_values)  
        probabilities = exp_q / np.sum(exp_q)  
        p_chosen = probabilities[choice1[t]]  
        # Update the log likelihood  
        LL += np.log(p_chosen) if p_chosen > 0 else -np.inf  
  
        # Learning at first level: update your eligibility trace according to  
        #  $e(a/S) = 1$  for the chosen action (a) in the_  
        ↪current state (S)  
        #  $e(a/S) = \lambda * e(a/S)$  for all other a, S pairs  
        E *= lambd # decay all the traces
```

```

E[S, choice1[t]] = 1

# Update prediction error without reward (because we are in the first
→level)
if model_based:
    # Implement SARSA update for model based learning
    # Keep in mind that choosing 1 at the first level (S1) leads to S2
→in approximately 70% of the trials
    # and choosing 2 at the first level (S1) leads to S3 in
→approximately 70% of the trials
    next_state = 1 if choice1[t] == 1 else 2 # Transition to S2 or S3
    next_action = choice2[t] # Next action taken
    PE = money[t] + Q[next_state, next_action] - Q[S, choice1[t]]

else:
    # Implement Q-learning update for model free learning
    PE = money[t] + np.max(Q[state[t]]) - Q[S, choice1[t]]

# update Q values according to  $Q = Q + \eta * \text{prediction error} * E$ 
→eligibility
Q += eta * PE * E

# Get your current state for the second level (S2 or S3)
S = state[t]

# Stop if trial was missed at the second level. Missed trials will have
→a value of -1
if choice2[t] == -1:
    continue

# Second level choice likelihood: compute likelihood of choice at the
→second state (S2 or S3).
# Your likelihood should be a softmax function.
q_values = Q[S]
exp_q = np.exp(beta * q_values)
probabilities = exp_q / np.sum(exp_q)
p_chosen = probabilities[choice2[t] - 1]
# Update your log likelihood
LL += np.log(p_chosen) if p_chosen > 0 else -np.inf

# Learning at second level: update your eligibility trace according to
#  $e(a/S) = 1$  for the chosen action (a) in the
→current state (S)
#  $e(a/S) = \lambda * e(a/S)$  for all other a, S pairs
E *= lambda
E[S, choice2[t]] = 1

```

```

        # Update the prediction error with reward because we are in the second
        ↪ level
        # NOTE: This update IS NOT dependent on the next state because we are
        ↪ in the final state
        PE = money[t] - Q[S, choice2[t] ]

        # update Q values according to  $Q = Q + \eta * \text{prediction error} * \text{eligibility}$ 
        ↪ eligibility
        Q += eta * PE * E

    return -LL

```

Question: Is the prediction error (δ) update in the second stage fundamentally similar or different between Q-learning and SARSA? Explain your answer.

[your answer here] The prediction error (δ) update in the second stage for both Q-learning and SARSA are fundamentally different because Q-learning considers the best possible future reward, while SARSA considers the expected reward according to the current policy's action.

Question: Which of these two algorithms is considered on-policy, which is off-policy, and why?

[your answer here] Q-learning is off-policy algorithm because it updates its predictions based on the maximum expected future rewards, regardless of the policy being followed. On the contrary, SARSA is an on-policy algorithm because it updates its predictions based on the action actually taken by the policy.

For each subject, load in their data as described at the beginning of the assignment (sub-0.npz to sub-4.npz). The sub-0.npz file contains sample data, while the rest are experimental data collected from other students at PNI. Use SciPy's `minimize` function (imported at the beginning of the problem set) to fit the two models to each of the subjects. You may also want to keep of the number of trials completed by each subject.

```

[11]: # Set some parameters
np.random.seed(1312)
params = [.5, 50, .5]
bounds = [(0, 1), (0, 100), (0, 1)]
sub_fns = ['sub-0.npz', 'sub-1.npz', 'sub-2.npz',
           'sub-3.npz', 'sub-4.npz']

# Example solver method for SciPy's minimize
method = 'TNC'

# Loop through subjects, load data, and fit models:

results = {}
for sub_fn in sub_fns:
    data = np.load(sub_fn)

```



```

c1 = data['choice1']
c2 = data['choice2']
s = data['state']
m = data['money']

if sub_fn not in results:
    results[sub_fn] = {}
    #Fit model Q-learning model
    #nll = lambda params: rl_nll(params, s, c1, c2, m)
    nll_q_learning = lambda params: rl_nll(*params, state=s, choice1=c1,
↪choice2=c2, money=m, model_based=False) #to unpack parmas and avoid error
    # Perform the optimization
    res_q_learning = minimize(nll_q_learning, params, bounds=bounds,
↪method=method)

    # Store the results
    results[sub_fn]['Q-learning'] = {'x': res_q_learning.x, 'fun':
↪res_q_learning.fun, 'n_trials': len(s)}

    #Fit model SARSA model
    nll_sarsa = lambda params: rl_nll(*params, state=s, choice1=c1, choice2=c2,
↪money=m, model_based=True)
    res_sarsa = minimize(nll_sarsa, params, bounds=bounds, method=method)
    results[sub_fn]['SARSA'] = {'x': res_sarsa.x, 'fun': res_sarsa.fun,
↪'n_trials': len(s)}

# Display the results
results

```

```

[11]: {'sub-0.npz': {'Q-learning': {'x': array([ 0.5, 50. ,  0.5]),
    'fun': 11917.841261874572,
    'n_trials': 500},
    'SARSA': {'x': array([ 0.5, 50. ,  0.5]),
    'fun': 12631.807684700607,
    'n_trials': 500}},
    'sub-1.npz': {'Q-learning': {'x': array([ 0.5, 50. ,  0.5]),
    'fun': 5148.594893050067,
    'n_trials': 201},
    'SARSA': {'x': array([ 0.5, 50. ,  0.5]),
    'fun': 5221.447325577614,
    'n_trials': 201}},
    'sub-2.npz': {'Q-learning': {'x': array([ 0.5, 50. ,  0.5]),
    'fun': 5325.6821979568085,
    'n_trials': 201},
    'SARSA': {'x': array([ 0.5, 50. ,  0.5]),
    'fun': 4993.45715367874,
    'n_trials': 201}},

```

```
'sub-3.npz': {'Q-learning': {'x': array([ 0.5, 50. ,  0.5]),
  'fun': 5190.685663720242,
  'n_trials': 201},
  'SARSA': {'x': array([ 0.5, 50. ,  0.5]),
  'fun': 5162.863956078575,
  'n_trials': 201}},
'sub-4.npz': {'Q-learning': {'x': array([ 0.5, 50. ,  0.5]),
  'fun': 4063.5999953850346,
  'n_trials': 201},
  'SARSA': {'x': array([ 0.5, 50. ,  0.5]),
  'fun': 4142.5298465596225,
  'n_trials': 201}}}
```

Use Bayesian information criterion (BIC) to compare which is the best-fitting model for each subject. Compute BICs using the following formula:

$$BIC = -2 * \log\text{-likelihood} + \ln(\text{number of trials}) * \text{number of parameters}$$

where $\ln()$ is the natural logarithm. BIC is defined here on the deviance scale, which means that lower values are better. **Question:** Which model fits each subject's behavior best?

```
[12]: # Compute BIC for each model and subject:
num_parameters = 3
bic_values = {}

# Iterate through the results to calculate BIC for each subject and model
for sub_fn, models in results.items():
    bic_values[sub_fn] = {}
    for model_name, model_results in models.items():
        log_likelihood = -model_results['fun']
        n_trials = model_results['n_trials']
        BIC = -2 * log_likelihood + np.log(n_trials) * num_parameters
        bic_values[sub_fn][model_name] = BIC

bic_values
```

```
[12]: {'sub-0.npz': {'Q-learning': 23854.32634804441, 'SARSA': 25282.25919369648},
  'sub-1.npz': {'Q-learning': 10313.099700824312, 'SARSA': 10458.804565879405},
  'sub-2.npz': {'Q-learning': 10667.274310637795, 'SARSA': 10002.824222081657},
  'sub-3.npz': {'Q-learning': 10397.281242164661, 'SARSA': 10341.637826881328},
  'sub-4.npz': {'Q-learning': 8143.109905494246, 'SARSA': 8300.969607843423}}
```

[your answer here] Lower BIC values indicate a better fit. For sub-0: Q-model fits best

For sub-1: Q-model fits best

For sub-2: SARSA fits best

sub-3: SARSA-model fits best

sub-4: Q-model fits best

1.0.2 Problem 2: Cliff walking

Consider the grid world shown below. This is a standard undiscounted, episodic task, with start and goal states, and the usual actions causing movement up, down, right, and left. Reward is -1 on all transitions except those into the region marked “The Cliff.” Stepping into this region incurs a reward of -100 and sends the agent instantly back to the start.

Two paths are marked: an optimal path which incurs the least costs on the way to the goal, and a roundabout (but safe) path that walks farthest from the cliff.

Question: Which algorithm, SARSA or Q -learning, would learn either path, and why?

[your answer here] SARSA takes the action selection policy into account even while learning, if the exploration policy has some chance of selecting a risky action near the cliff, it might receive the -100 penalty by occasionally falling off the cliff. Therefore, SARSA is more likely to learn the safer, roundabout path.

The Q -model learns the optimal policy. That means it learns the value of the best possible action to take, regardless of the exploration policy being used. Even if the exploration steps occasionally lead to the cliff, the Q -values will still converge to those representing the optimal policy, which in this case is the shortest path to the goal. It always updates its values towards the maximum future reward, irrespective of the exploration missteps.

Question: When behaving according to the softmax of the learned Q values, which path would an agent prefer? (Consider the parameter β and the stability of the environment.)

[your answer here] If beta is low, the action selection is more exploratory, as the differences in the Q values lead to smaller differences in action probabilities. The agent’s choices will be more random, and the probability of selecting any given action is less sensitive to the Q value differences. In this case, the agent may sometimes choose the safer path simply because the action probabilities are more uniform. If beta is high, the action selection is more exploitative. The agent is more likely to choose actions with the highest Q values. If the environment is stable and the Q values have converged correctly, the agent will prefer the path that leads to the highest expected reward, which would likely be the optimal path along the edge of the cliff.

Question: Can you explain why on-policy methods might be superior for learning real-world motor behavior?

[your answer here] On-policy methods are continuously updated based on the actual consequences of actions taken, which is vital for motor tasks where feedback from the physical environment is a key component of learning. This also helps in adapting to unpredictable environment. Additionally, on-policy methods can integrate complexity into a task.

References

- Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011). Model-based influences on humans’ choices and striatal prediction errors. *Neuron*, 69(6), 1204–1215. <https://doi.org/10.1016/j.neuron.2011.02.027>
- Sutton, R. S., & Barto, A. G. (1992). Reinforcement Learning: An Introduction. MIT Press.