# Homework #3
# **Winter is coming...**

## Problem

### Description

For this assignment, you will build a Sarsa agent which will learn policies in the Frozen Lake environment. You will be given randomized Frozen Lake maps, with corresponding sets of parameters to train your Sarsa agent; if your agent is implemented correctly, after training for the specified number of episodes it will produce the same policies (which are not necessarily optimal policies) as the grader.

| | | | |
|---|---|---|---|
| S | F | F | F |
| F | H | F | H |
| F | F | F | H |
| H | F | F | G |

Example Frozen Lake Map

[Frozen Lake](#) is a grid world environment that is highly stochastic, where the agent must cross a slippery frozen lake which has deadly holes to fall through. The agent begins in the starting state (S) and is given a reward of 1 if it reaches the goal state (G). The agent can take one of four possible moves at each state (left, down, right, or up). The frozen cells (F) are slippery, so the agent's actions succeed only ⅓ of the times, while the other ⅔ are split evenly in orthogonal directions. If the agent lands in a hole (H), then the episode terminates, and the agent is given a reward of 0.

## Sarsa $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$

Sarsa is a model-free on-policy reinforcement learning algorithm that solves the control problem through trial-and-error learning. It is model-free because, unlike with value iteration and policy iteration, it doesn't need or use an MDP. It is on-policy because it learns about the same policy that generates behaviors. The algorithm estimates the action-value function $Q_\pi$ of the behavior policy $\pi$, and uses an exploration strategy to improve $\pi$ while increasing the policy's greediness.

### Procedure

1. Install OpenAI Gym
2. Implement a Sarsa agent
3. Train your agent given the specified parameters
4. Export the policy and enter it into the RLDM site

## Notes

- You must use Python, NumPy, and OpenAI Gym for this homework
- Use the provided seed for both Gym and NumPy
- Initialize the agent's Q-table to zeros
- To avoid any unexpected behavior, setup the Gym environment with `gym.envs.toy_text.frozen_lake.FrozenLakeEnv().unwrapped`
- To set up the environment with a custom map, use the `desc` variable
- You must train your agent with an epsilon-greedy exploration strategy, using NumPy's `numpy.random.randint` function to select random actions

# Examples

The following examples can be used to verify that your agent is implemented correctly.

- Input: `amap=SFFFHFFFFFFFFFFG, gamma=1.0, alpha=0.25, epsilon=0.29, n_episodes=14697, seed=741684,` Output: `^,v,v,>,<,>,>,v,v,v,>,v,>,>,>,<`
- Input: `amap=SFFFFHFFFFFFFFFFFFFFFFFFFG, gamma=0.91, alpha=0.12, epsilon=0.13, n_episodes=42271, seed=983459,` Output: `^,>,>,>,>,<,>,>,>,v,v,v,>,>,v,v,>,>,>,>,v,>,>,^,<`
- Input: `amap=SFFG, gamma=1.0, alpha=0.24, epsilon=0.09, n_episodes=49553, seed=202404,` Output: `<,<,v,<`

# Resources

The concepts explored in this homework are covered by:

- Lectures
    - Lesson 4: Convergence
- Readings
    - Sutton 2018, Chapter 6 (6.4 Sarsa: On-policy TD Control)

# Submission Details

**The due date is indicated on the Canvas page for this assignment.**

Make sure you have set your timezone in Canvas to ensure the deadline is accurate.

To complete the assignment calculate answers to the specific problems given and submit results at https://rldm.herokuapp.com

∎ ∎ ∎