

Homework

- Use the CliffWalking domain from OpenAI gym
 - See Example 6.6, pg 132 in Sutton and Barto [2018]
- Modify the TD(λ) algorithm presented to implement SARSA(λ)
 - The only difference here is that there is an eligibility trace for each **state-action** pair!
 - See the **first edition of Sutton and Barto** for more info
 - Use ε -greedy policies with $\varepsilon = 0.1$ and a learning rate of $\alpha = 0.5$
 - Run SARSA(λ) on the domain for $\lambda = \{0, 0.3, 0.5, 0.7, 0.9\}$ for 500 episodes
 - Record the current estimate of the Q-value function after each episode

By next week's lecture, submit on Moodle:

1. Perform a single run of the algorithm. After each episode plot the value function (take $\max_a Q(s, a)$) learned so far as a heatmap for each λ side by side. Ensure the visualisation aligns with the layout of the domain. This should result in 500 separate plots/images. Turn these images into an animation/video and submit it.
2. Your code