



Aprendizaje Reforzado

Maestría en Ciencia de Datos, DC - UBA

Julián Martínez
Javier Kreiner



Diferencias Temporales (programación)

- Predicción TD
- Sarsa
- Q-learning



Problema (programación)

- Implementar expected-Sarsa, Sutton sección 6.6, es igual Q-learning, pero la ecuación de update es, usarlo para resolver Windy Gridworld y Cliff Environment :

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \sum \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right], \end{aligned}$$