

V023723245 Lezama Luis

V025793252 Ramírez Coalbert

# Análisis del método **Value Iteration** Vs. **Policy Iteration**

$\gamma$	Value Iteration Avg. Reward	Policy Iteration Avg. Reward	Value Iteration $V(s)$				Policy Iteration $V(s)$				Value Iteration $\pi(s)$	Policy Iteration $\pi(s)$
0.6	0.455	0.392	0.001348 0.003564 0.012909 0.	0.001830 0. 0.048072 0.123930	0.00597256 0.021899 0.10352 0.44764	0.001990 0. 0. 0.	0.001254 0.003348 0.011814 0.	0.001680 0. 0.045936 0.115923	0.005802 0.021056 0.101250 0.430138	0.001772 0. 0. 0.	1. 3. 2. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 1. 0.	2. 3. 2. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 2. 0.
0.7	0.453	0.392	0.004503 0.009610 0.027074 0.	0.005187 0. 0.079349 0.172448	0.012541 0.035720 0.140545 0.487266	0.00548 0. 0. 0.	0.004035 0.008719 0.023738 0.	0.004535 0. 0.074119 0.156069	0.011894 0.033750 0.135266 0.460865	0.004602 0. 0. 0.	1. 3. 2. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 1. 0.	1. 3. 2. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 2. 0.
0.8	0.455	0.392	0.015434 0.026853 0.058413 0.	0.015590 0. 0.133783 0.246537	0.027440 0.059780 0.196735 0.544195	0.015680 0. 0. 0.	0.013074 0.023083 0.048186 0.	0.0126927 0. 0.120622 0.212973	0.024916 0.054939 0.184414 0.503864	0.012029 0. 0. 0.	2. 3. 2. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 1. 0.	2. 3. 2. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 2. 0.
0.9	0.782	0.677	0.068890 0.091854 0.145436 0.	0.061414 0. 0.247496 0.379935	0.074409 0.112208 0.299617 0.639020	0.055807 0. 0. 0.	0.053108 0.071884 0.108574 0.	0.044974 0. 0.209059 0.304876	0.061798 0.097510 0.267153 0.572129	0.037020 0. 0. 0.	0. 3. 0. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 1. 0.	0. 3. 0. 3. 0. 0. 0. 0. 3. 1. 0. 0. 0. 2. 2. 0.

Se realizó la ejecución de los algoritmos de **Value Iteration** y **Policy Iteration** 10,000 veces cada uno para obtener las recompensas promedio, de lo cual se puede concluir que con un factor de descuento  $\gamma$  igual a 0.9 genera una mayor recompensa promedio, **0.782** con **Value Iteration** y **0.667** con **Policy Iteration** para la simulación del frozen lake, lo que permite encontrar con mayor eficiencia el estado final, sin embargo, para valores de  $\gamma$  menores a 0.8, se puede observar que la eficiencia disminuye casi a la mitad para ambos métodos, pero a partir de allí mantienen la misma recompensa promedio, donde se pueden apreciar cambios es en la tabla de valores  $V(s)$ , donde se puede notar que a medida que incrementa el valor del factor de descuento  $\gamma$ , también aumentan las valoraciones en la tabla