

State space S : 10 space-subdivision,
2D grid $\Rightarrow n_{\text{states}} = \text{space-subdivision}^2 =$
100

Action space A : $A = \{ \text{left, right, down, up, stay} \}$

Hyperparameters α, γ scenarios:

(1) Not improving/learning:

$$\alpha = (0; 0.1) ; \gamma = (0; 1]$$

(2) High variance but fast learning:

$$\alpha = [0.7; 1] ; \gamma = (0; 1]$$

(3) low variance and high long-term return:

$$\alpha = [0.1; 0.1] ; \gamma = [0.7; 1]$$

(4) High variance and high long-term return:

$$\alpha = [0.7; 1] ; \gamma = [0.7; 1]$$

long-term return of the optimal policy in 32.3.4ml:

Reward for golden fish: 50.

Reward for step: -2.

Reward Jelly: -10.

$$\Rightarrow R_{opt} = 50 - 2 \cdot \text{shortest path} \\ = 50 - 2 \cdot 10 = \underline{30!}$$

