





```
In [20]: stats = { "rewards": [[], []] }

# Note that the reward position is adjusted -- we added walls to the
# outer edges of the maze.
goal = (6, 6)
goal_state = goal[0] + maze.shape[1] + goal[1]

np.random.seed(42)
max_int = np.iinfo(np.int32).max

for i in range(20):
    # Using original SR from the random walk policy
    M_clamped, V_clamped, earned_rewards_clamped, traj_clamped, _, _ = actor_critic(
        analytical_state_repr,
        n_steps = 300,
        alpha = 0.05,
        gamma = 0.99,
        n_episodes = 1000,
        update_sr = False,
        start_func = normal_start,
        seed = np.random.randint(1, max_int)
    )
    stats["rewards"][0].append(earned_rewards_clamped)

    # Using re-learned SR (tuned towards finding the reward at (1, 1))
    M_relearned, V_relearned, earned_rewards_relearned, traj_relearned, learned_state_repr_relearned,
    n_steps = 300,
    alpha = 0.05,
    gamma = 0.99,
    n_episodes = 1000,
    update_sr = True,
    start_func = normal_start,
    seed = np.random.randint(1, max_int)
    )
    stats["rewards"][1].append(earned_rewards_relearned)

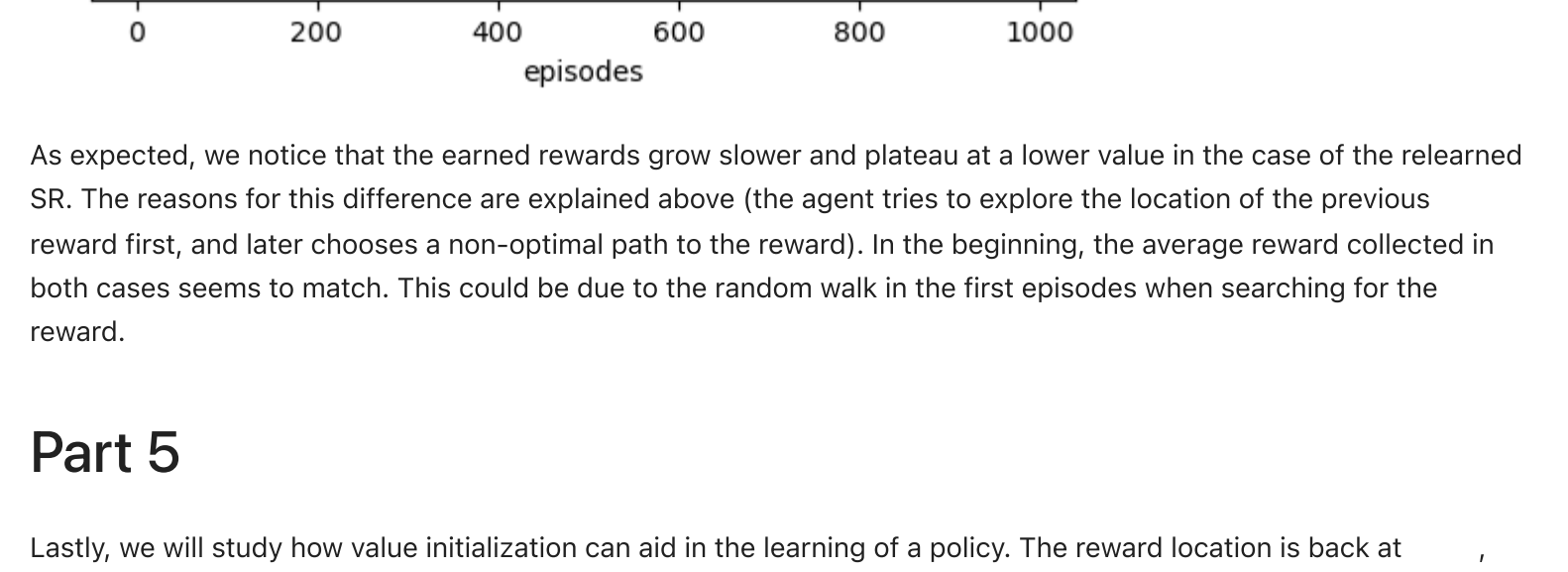
100% 1000/1000 (00:00-00:00, 1106.60it/s)
100% 1000/1000 (00:01-00:00, 618.74it/s)
100% 1000/1000 (00:00-00:00, 1086.09it/s)
100% 1000/1000 (00:01-00:00, 614.85it/s)
100% 1000/1000 (00:00-00:00, 1074.72it/s)
100% 1000/1000 (00:01-00:00, 535.66it/s)
100% 1000/1000 (00:00-00:00, 1262.61it/s)
100% 1000/1000 (00:01-00:00, 597.35it/s)
100% 1000/1000 (00:01-00:00, 843.22it/s)
100% 1000/1000 (00:01-00:00, 613.13it/s)
100% 1000/1000 (00:00-00:00, 1247.96it/s)
100% 1000/1000 (00:01-00:00, 567.77it/s)
100% 1000/1000 (00:00-00:00, 1190.11it/s)
100% 1000/1000 (00:01-00:00, 625.78it/s)
100% 1000/1000 (00:00-00:00, 1153.68it/s)
100% 1000/1000 (00:01-00:00, 509.90it/s)
100% 1000/1000 (00:00-00:00, 1107.00it/s)
100% 1000/1000 (00:01-00:00, 947.83it/s)
100% 1000/1000 (00:01-00:00, 550.73it/s)
100% 1000/1000 (00:00-00:00, 1241.93it/s)
100% 1000/1000 (00:01-00:00, 587.80it/s)
100% 1000/1000 (00:00-00:00, 1053.00it/s)
100% 1000/1000 (00:01-00:00, 573.41it/s)
100% 1000/1000 (00:00-00:00, 1060.69it/s)
100% 1000/1000 (00:01-00:00, 594.84it/s)
100% 1000/1000 (00:00-00:00, 1080.87it/s)
100% 1000/1000 (00:01-00:00, 575.56it/s)
100% 1000/1000 (00:00-00:00, 1045.34it/s)
100% 1000/1000 (00:01-00:00, 596.85it/s)
100% 1000/1000 (00:00-00:00, 1136.69it/s)
100% 1000/1000 (00:01-00:00, 563.86it/s)
100% 1000/1000 (00:00-00:00, 1258.96it/s)
100% 1000/1000 (00:01-00:00, 566.76it/s)
100% 1000/1000 (00:00-00:00, 1080.56it/s)
100% 1000/1000 (00:01-00:00, 572.24it/s)
100% 1000/1000 (00:00-00:00, 1314.15it/s)
100% 1000/1000 (00:01-00:00, 606.93it/s)
100% 1000/1000 (00:00-00:00, 1133.95it/s)
100% 1000/1000 (00:01-00:00, 620.88it/s)
```

```
In [21]: fig, ax = plt.subplots(ncols = 2, figsize = (12, 6))

plot_path(maze, int2state(traj_clamped[-1]), start = start, end = goal, ax = ax[0])
lm0 = ax[0].imshow(
    (analytical_state_repr @ V_clamped).reshape(maze.shape),
    cmap = "hot",
    alpha = 0.8,
    vmin = 0,
    vmax = 10
)
fig.colorbar(lm0, shrink = 0.6)

plot_path(maze, int2state(traj_relearned[-1]), start = start, end = goal, ax = ax[1])
lm1 = ax[1].imshow(
    (learned_state_repr_relearned @ V_relearned).reshape(maze.shape),
    cmap = "hot",
    alpha = 0.8,
    vmin = 0,
    vmax = 10
)
ax[1].set_title("relearned")
fig.colorbar(lm1, shrink = 0.6)

plt.show()
```



In the 1st case (using the original SR from the random walk policy), the agent does not know where the reward is located and proceeds to explore the environment and locate the reward position as usual. In the 2nd case (using the re-learned SR that is tuned towards finding the reward at (1, 1)), the agent first tries to search for the reward at its usual location, and then proceeds to explore the environment from there when the reward is not found. So, in the 2nd case, we see the agent exploring the environment where the reward used to be (indicated by the bright colors in the upper left side of the maze). The final trajectory in this case is not optimal (as the exploration of the optimal path is not encouraged), and it branches off the path to the previous goal location.

```
In [22]: reward_stats = np.stack(stats["rewards"])

mean_base = reward_stats.mean(axis = 1)
mean = np.stack([
    np.arange(len(mean[0])),
    running_mean(mean_base[0], 10),
    running_mean(mean_base[1], 10),
])

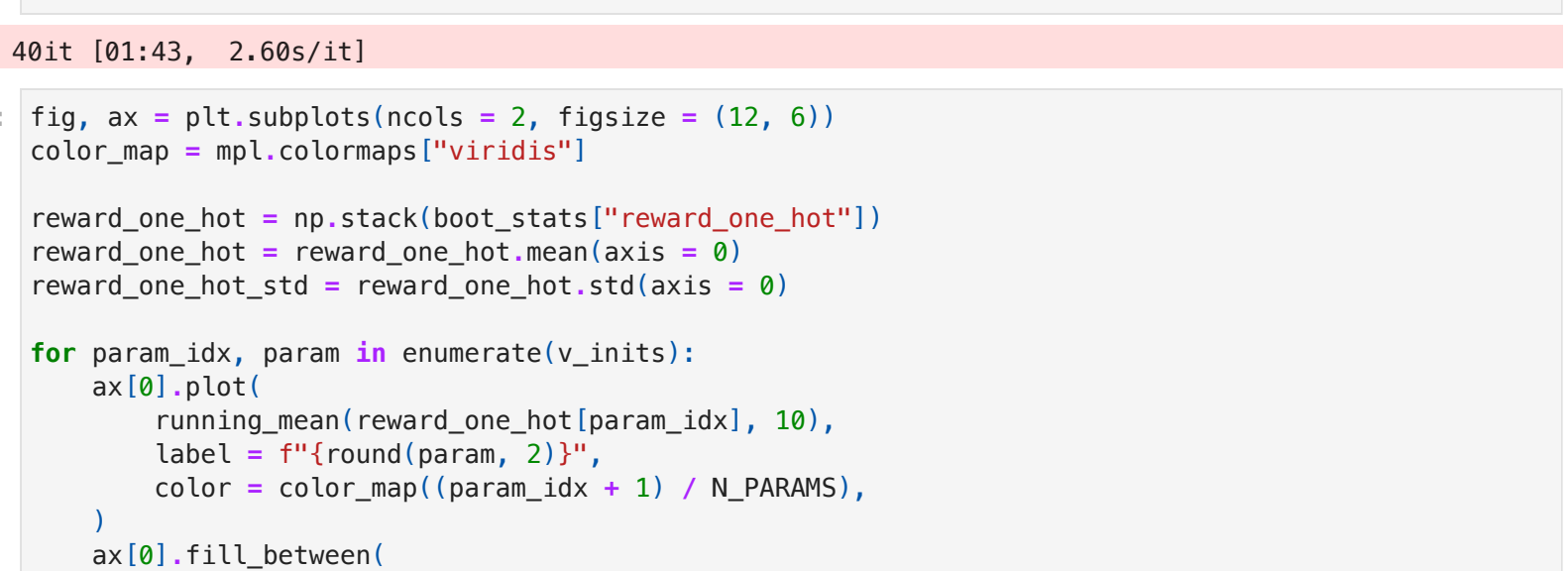
std_base = reward_stats.std(axis = 1)
std = np.stack([
    np.arange(len(std_base[0], 10),
    running_mean(std_base[0], 10),
    running_mean(std_base[1], 10),
])

tab10 = mpl.colormaps["tab10"]

plt.plot(mean[0], label = "clamped")
plt.fill_between(
    np.arange(len(mean[0])),
    mean[0] - std[0],
    mean[0] + std[0],
    color = tab10(0),
    alpha = 0.2,
)

plt.plot(mean[1], label = "relearned")
plt.fill_between(
    np.arange(len(mean[1])),
    mean[1] - std[1],
    mean[1] + std[1],
    color = tab10(1),
    alpha = 0.2,
)

plt.xlabel("episodes")
plt.ylabel("earned rewards")
plt.legend(loc = "center left", bbox_to_anchor = (1, 0.5))
plt.show()
```



As expected, we notice that the earned rewards grow slower and plateau at a lower value in the case of the relearned SR. The reasons for this difference are explained above (the agent tries to explore the location of the previous reward first, and later chooses a non-optimal path to the reward). In the beginning, the average reward collected in both cases seems to match. This could be due to the random walk in the first episodes when searching for the reward.

## Part 5

Lastly, we will study how value initialization can aid in the learning of a policy. The reward location is back at (1, 1). We always start at the original starting position, and use the SR from a random-walk policy as our representation. So far, we have initialized our weights  $w$  with  $\mathcal{U}(-1, 1)$ . Experiment with different initializations along with both the 1-hot representation and the SR. Try a couple of extreme points (like 4-5 different values) from 0 to 90 as your initialization. What do you observe, why do you think some values help while others hurt?

```
In [23]: # Reset the goal location (again, adjusted by our inclusion of outer walls in the maze)
goal = (2, 2)
goal_state = state2int(goal)

In [24]: # Define parameters to run multiple models
N_BOOT = 10
N_EPISODES = 1000
N_STEPS = 300
ALPHA = 0.05
GAMMA = 0.99

# Generate initial weight values in logspace
v_inits = np.logspace(0, np.log10(90 + 1), num = N_PARAMS, base = 10)
v_inits -= 1 # Adjust for log

np.random.seed(42)
max_int = np.iinfo(np.int32).max

In [25]: # "reward_one_hot": np.zeros(N_BOOT, N_PARAMS, N_EPISODES), "reward_sr": np.zeros(N_BOOT,
N_PARAMS, N_EPISODES)

for v_init_idx, boot_idx in tqdm(product(range(N_PARAMS), range(N_BOOT))):
    M_one_hot, V_one_hot, earned_rewards_one_hot, traj_one_hot, _, _ = actor_critic(
        analytical_state_repr,
        n_steps = N_STEPS,
        alpha = ALPHA,
        gamma = GAMMA,
        n_episodes = N_EPISODES,
        v_init = v_inits[v_init_idx],
        start_func = normal_start,
        quiet=True,
        seed = np.random.randint(1, max_int)
    )
    boot_stats["reward_one_hot"][boot_idx, v_init_idx] = earned_rewards_one_hot

40it [02:41, 4.03s/it]
```

```
In [26]: # Run multiple fits for random-walk policy state representation
for v_init_idx, boot_idx in tqdm(product(range(N_PARAMS), range(N_BOOT))):
    M_sr, V_sr, earned_rewards_sr, traj_sr, _, _ = actor_critic(
        analytical_state_repr,
        n_steps = N_STEPS,
        alpha = ALPHA,
        gamma = GAMMA,
        n_episodes = N_EPISODES,
        v_init = v_inits[v_init_idx],
        start_func = normal_start,
        quiet=True,
        seed = np.random.randint(1, max_int)
    )
    boot_stats["reward_sr"][boot_idx, v_init_idx] = earned_rewards_sr

40it [01:43, 2.60s/it]
```

```
In [27]: fig, ax = plt.subplots(ncols = 2, figsize = (12, 6))
color_map = mpl.colormaps["viridis"]

reward_one_hot = np.stack(boot_stats["reward_one_hot"])
reward_one_hot = reward_one_hot.mean(axis = 0)
reward_one_hot_std = reward_one_hot.std(axis = 0)

for param_idx, param in enumerate(v_inits):
    ax[0].plot(
        running_mean(reward_one_hot[param_idx], 10),
        label = f"({round(param, 2)})",
        color = color_map((param_idx + 1) / N_PARAMS),
    )
    ax[0].fill_between(
        np.arange(len(reward_one_hot[param_idx]) - 10 + 1),
        running_mean(reward_one_hot[param_idx] - reward_one_hot_std[param_idx], 10),
        running_mean(reward_one_hot[param_idx] + reward_one_hot_std[param_idx], 10),
        color = color_map(param_idx / N_PARAMS),
        alpha = 0.1
    )

ax[0].set_ylim(-1, 10)
ax[0].set_title("One-hot state representation")
ax[0].set_xlabel("epochs")
ax[0].set_ylabel("earned reward")

reward_sr = np.stack(boot_stats["reward_sr"])
reward_sr = reward_sr.mean(axis = 0)
reward_sr_std = reward_sr.std(axis = 0)

for param_idx, param in enumerate(v_inits):
    ax[1].plot(
        running_mean(reward_sr[param_idx], 10),
        label = f"({round(param, 2)})",
        color = color_map((param_idx + 0) / N_PARAMS),
    )
    ax[1].fill_between(
        np.arange(len(reward_sr[param_idx]) - 10 + 1),
        running_mean(reward_sr[param_idx] - reward_sr_std[param_idx], 10),
        running_mean(reward_sr[param_idx] + reward_sr_std[param_idx], 10),
        color = color_map(param_idx / N_PARAMS),
        alpha = 0.1
    )

ax[1].set_ylim(-1, 10)
ax[1].set_title("Random-walk policy state representation")
ax[1].set_xlabel("epochs")
ax[1].set_ylabel("earned reward")

plt.legend(loc = "center left", bbox_to_anchor = (1, 0.5), title = f"$V_{init}$")
plt.show()
```



The most noticeable difference between the different initial weight values is the speed of the convergence. The initial weights act as the agent's prior knowledge. They define the starting estimate of the value of each state before any learning occurs. If these estimates are far from the true values, they can hinder learning. High initial values introduce a bias that overestimates the agent's future rewards, and the agent may start prioritizing states based on inflated values rather than actual rewards. In general, when the initial values are too high, the agent wastes many episodes reducing these inflated estimates, delaying convergence. Moreover, it is more difficult to differentiate between the good and the bad states, which can hinder exploration, as the agent might stick to suboptimal paths due to misleading initial values. Smaller or zero initial values are better, providing a neutral or slightly optimistic starting point that aligns with typical rewards.

In the case of one-hot encoding state representation, the algorithm does not converge for the highest explored initial values. This is the case for the 1000 episodes used for this analysis; that is, the algorithm would converge with more iterations.

As in the previous cases, there is a noticeable convergence delay between the one-hot state representation and the random-walk policy state representation (the reasons for this were explained previously).

```
In [28]:
```