

Evaluation of Optimal Decision Making with Dead-ends in High-Risk Environments

Alex Arron Kashi and Marcel Torne Villasevil

Advised By

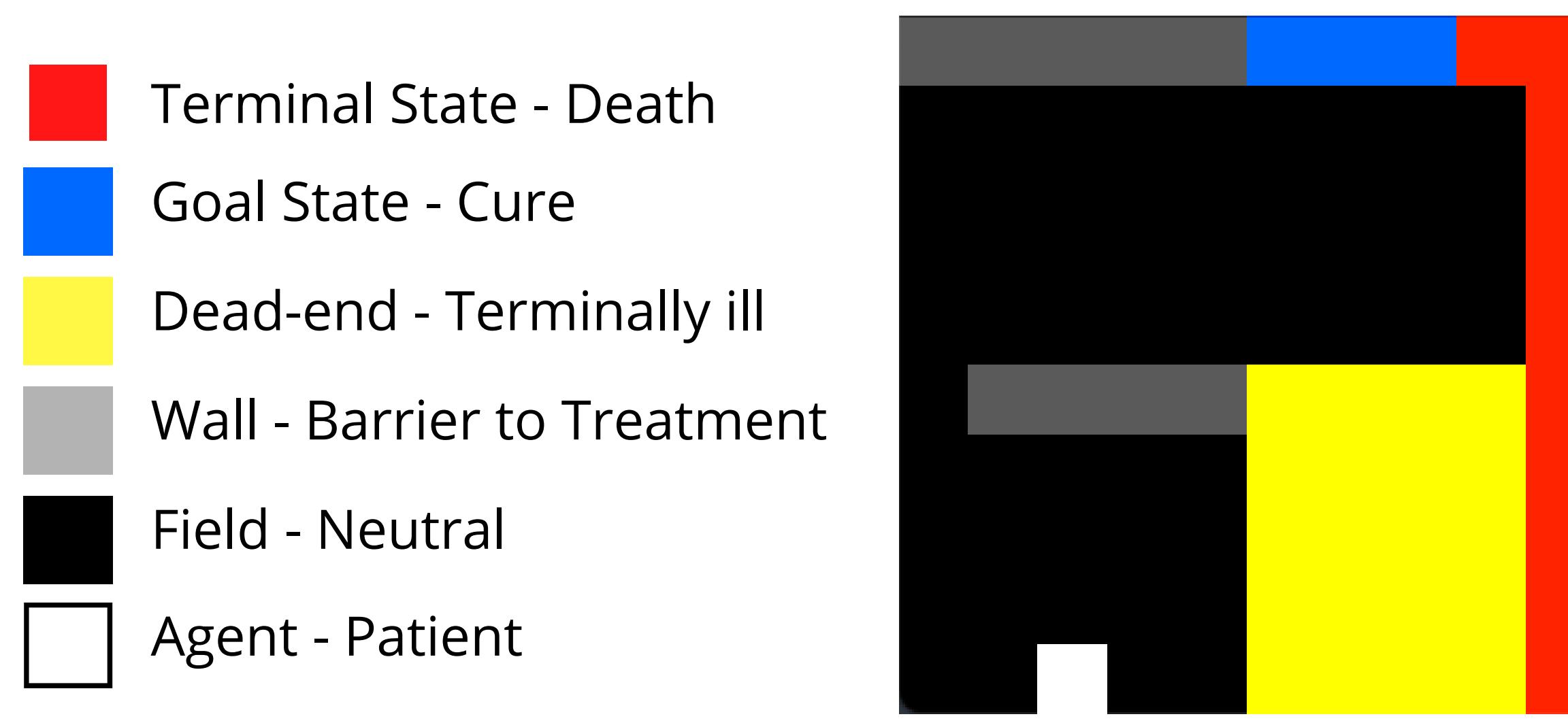
Professor Susan Murphy PhD., Harvard
Raaz Dwivedi PhD., Harvard
Eura Shin, Harvard

Motivation

Dead-ends are **states** that will certainly **lead** the agent to a **negative terminal state**. Dead-end identification is a significant problem that occurs in a variety of fields such as **healthcare** (with identifying these in treatments), **robotics** (identifying crashes before these happen), **stock markets** (identifying when the stocks are just going to lose value), and many more. Finally, it is a non-trivial problem, and in this project, we study it further.

Environment

Adapted from Fatemi et al.¹.



State Space

Agent position (y, x)

Action Space

No-op, Up, Down, Left, Right

Transitions

Death drag - 20% chance of a forceful right

In dead-end? - 70% chance of a forceful right, otherwise no-op

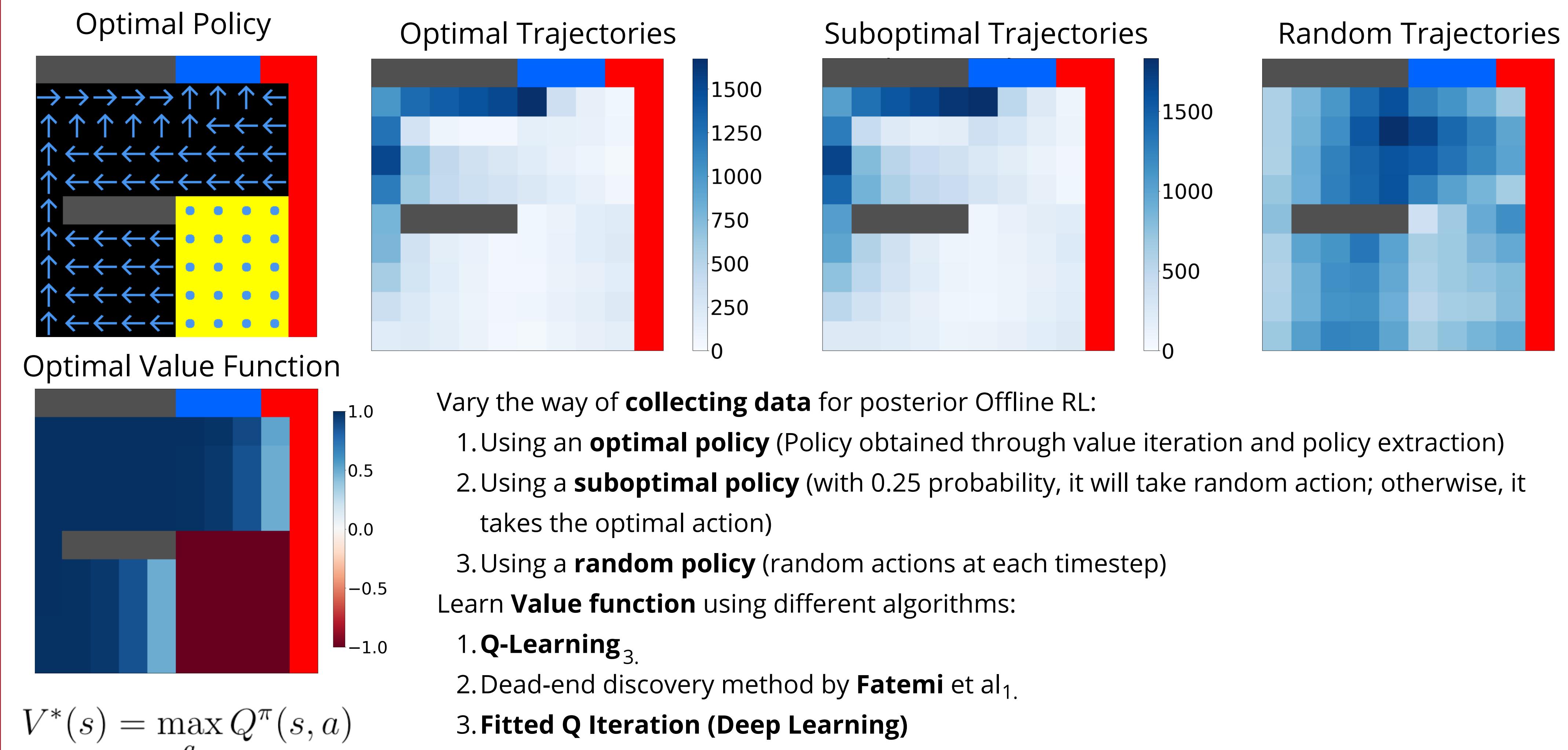
Rewards

Goal state (1)

Terminal State (-1)

Objective

- Identify dead-ends (propose different method)
- Evaluate how the data collection policy affects each method
- Learn a policy that can drive agents to recovery states



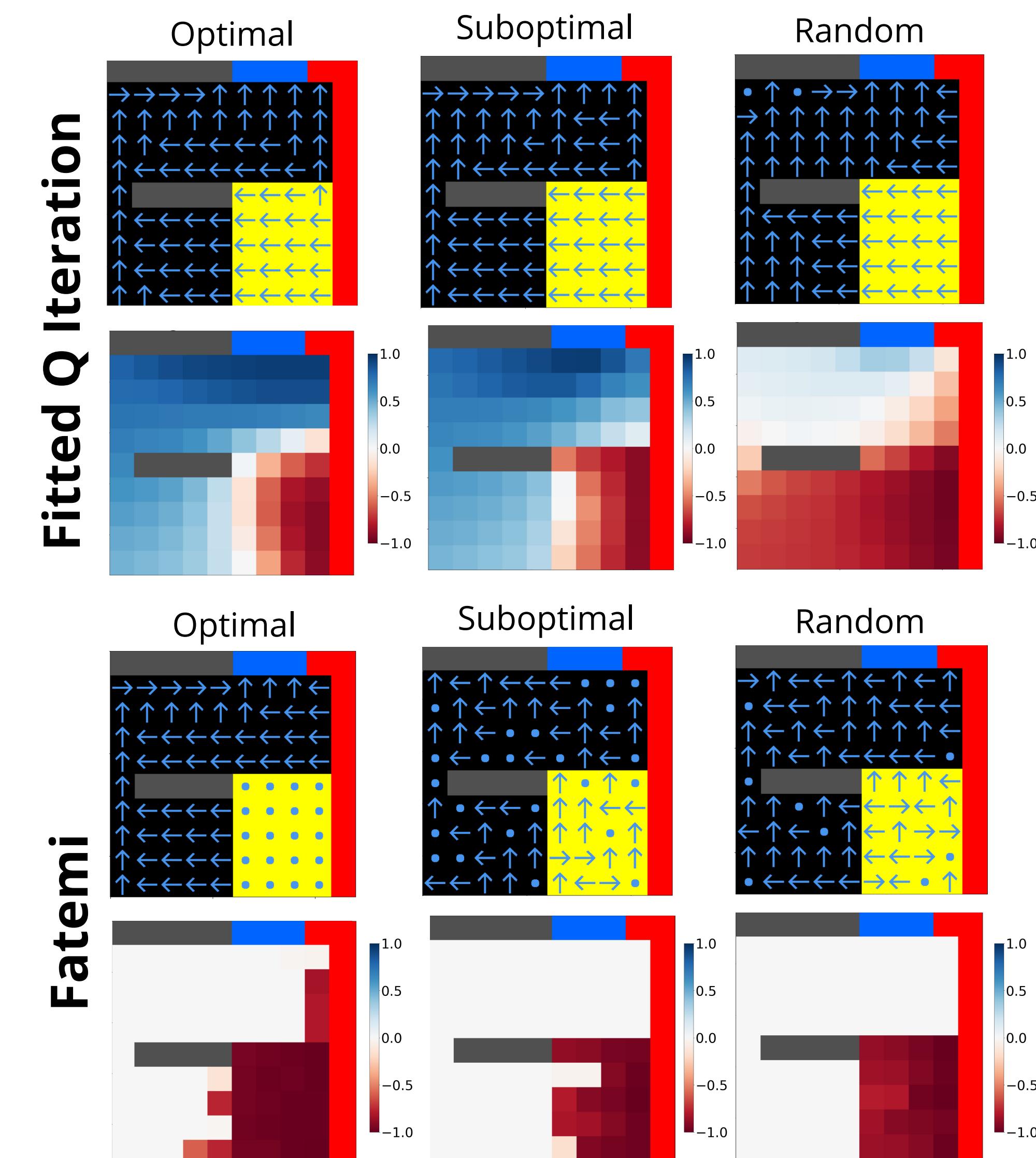
Results

Experimental Setup:

- Collect 2500 trajectories
- Max length of 100 steps.

Analysis:

- Fatemi's method produces results that are worse than Q-learning. We can deduce that knowing which states lead to recovery gives valuable information about dead ends.
- Regarding Q-Learning, no matter how the data was generated, it still recovers a reasonable estimation of the value function.
- Regarding Fitted Q Iteration, this is not the case. The best approximation is found for the data collected using the optimal policy.
- Concerning Fatemi's method, all recover a proper value function, but the best is found for the random data.



Conclusion

All methods studied in this project, Fitted Q Iteration, Q-Learning, and the method proposed by Fatemi et al., all manage to identify the dead-ends in the proposed toy environment. Furthermore, we saw that with the different ways that the data could be collected, most of these methods managed to recover the Q-values. Nevertheless, we observe that the results are the most consistent for Q-Learning. In addition, this method needs less training time than Fitted Q Iteration, and for this reason, given out experiments, we would propose Q-Learning for the dead-end identification.

Methods

Algorithm 14 Fitted Q iteration

```

1: procedure FQI
2:   Initialize the neural network weights  $\phi$  randomly
3:   while not terminate do
4:     Collect dataset  $\mathcal{D}$  using an arbitrary policy to interact with the environment
5:     for  $i = 1 \dots$  do
6:       Compute the target  $y \leftarrow r(s, a) + \gamma \max_{a'} Q_\phi(s', a') \forall (s, a, r, s') \in \mathcal{D}$ 
7:       Update the network  $\phi \leftarrow \arg \min_{\phi} \frac{1}{|\mathcal{D}|} \sum_{(s, a, r, s') \in \mathcal{D}} \|y - Q_\phi(s, a)\|^2$ 
8:     end for
9:   end while
10: end procedure

```

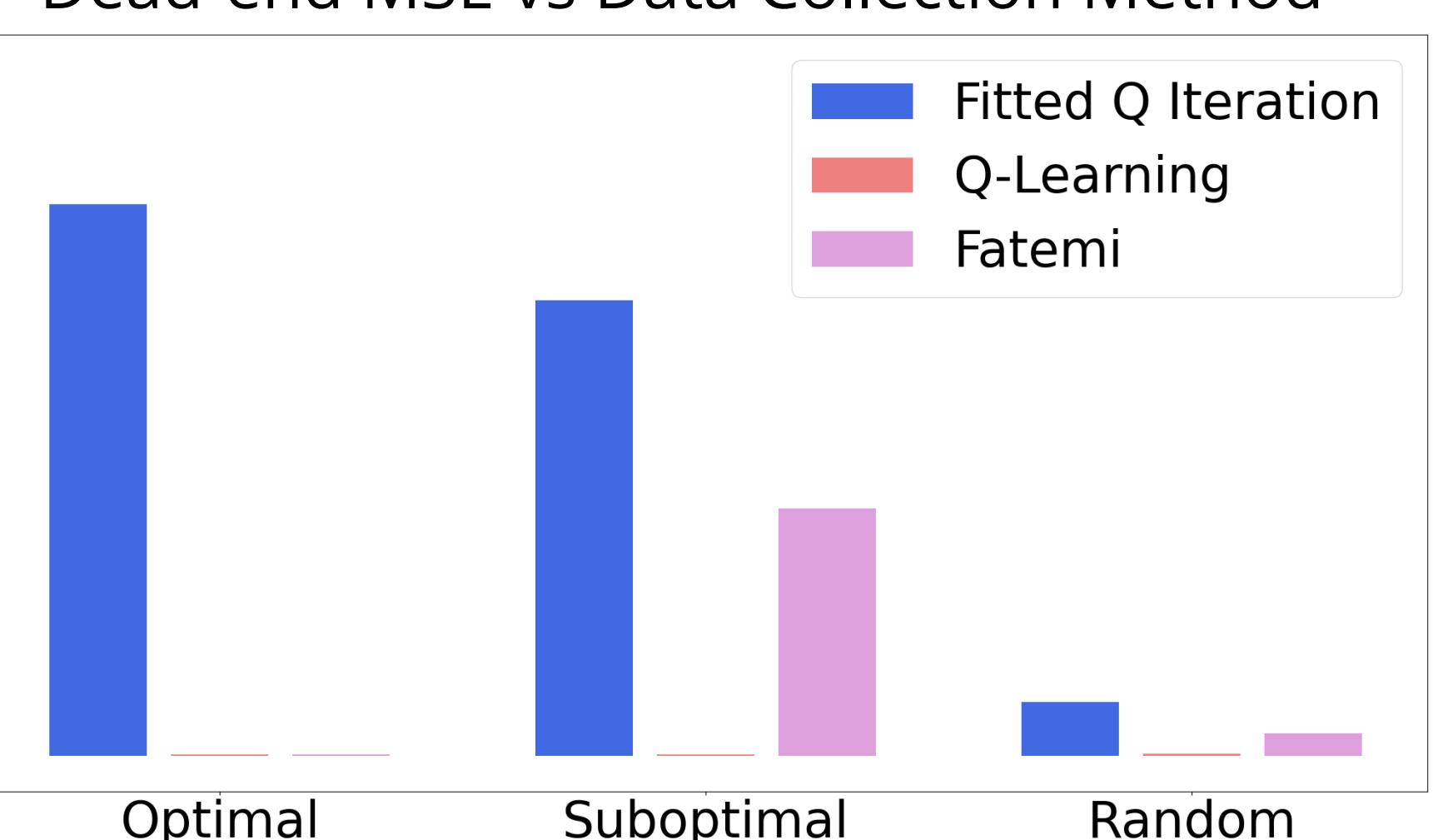
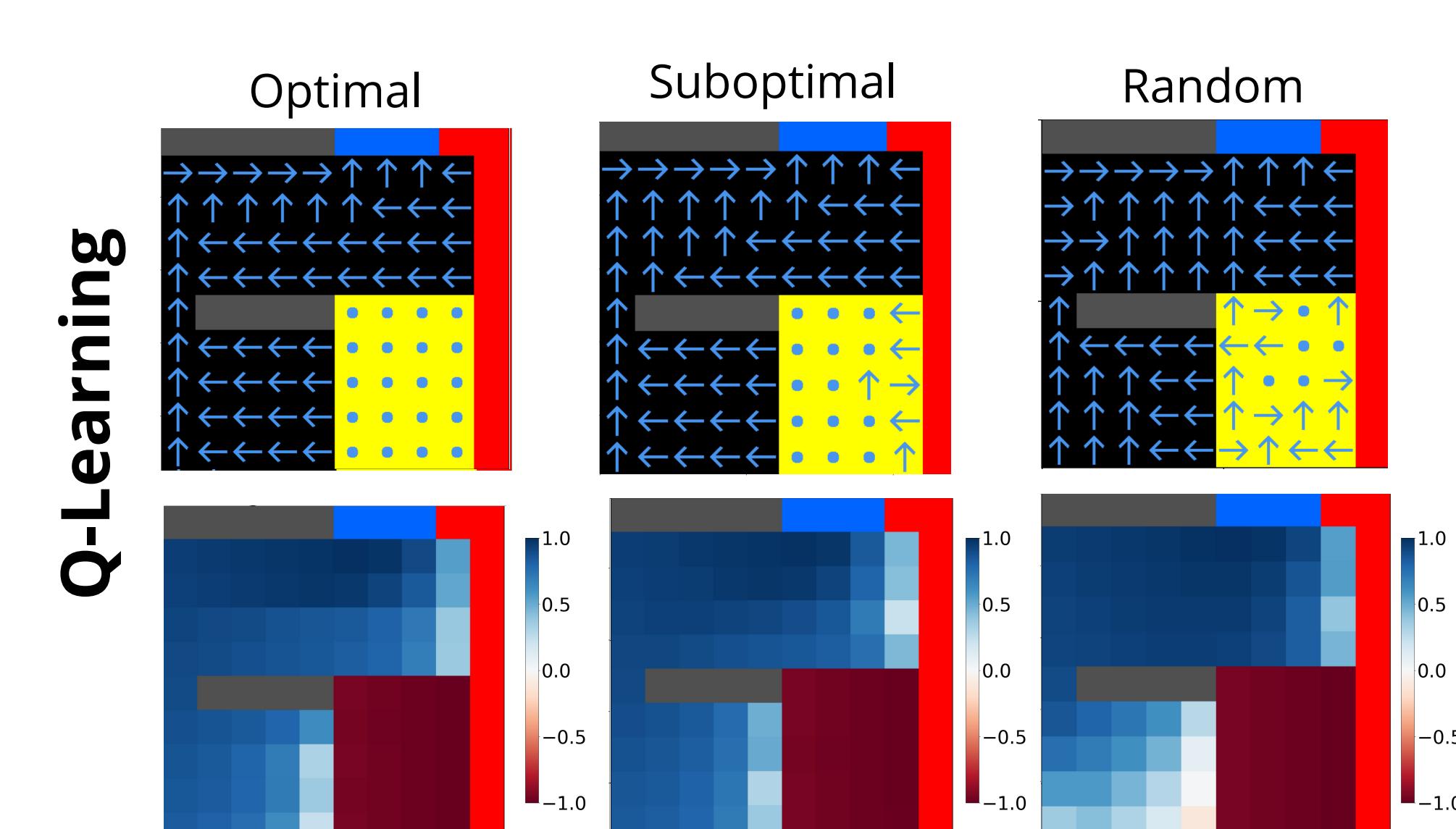
Algorithm 12 Value Iteration

```

1: procedure VALUEITERATION( $\pi, \epsilon$ )
2:   Initialize  $V_1(s) = 0 \forall s \in \mathcal{S}$ 
3:   for  $i = 1 \dots$  do
4:      $V_{i+1}(s) \leftarrow \max_a \sum_{s' \in \mathcal{S}} \mathcal{T}(s'|s, a)[r(s, a) + \gamma V_i(s')] \forall s \in \mathcal{S}$ 
5:     if  $\|V_{i+1} - V_i\|_\infty \leq \epsilon$  then
6:       break
7:     end if
8:   end for
9:   Set policy  $\pi(s) = \arg \max_a \mathbb{E}_\pi[r(s_t, a_t) + \gamma V^\pi(s_{t+1})|s_t = s, a_t = a] \forall s \in \mathcal{S}$ 
10: end procedure

```

Algorithms referenced from MIT 6.484 lecture notes².



References

- Fatemi, Mehdöli, et al. "Medical Dead-ends and Learning to Identify High-risk States and Treatments." Advances in Neural Information Processing Systems 34 (2021).
- Zhang-Wei Hong, Pulkit Agrawal, Lecture Notes of 6.484 Computational Sensorimotor Learning, MIT, 2022
- Sutton, Richard S., and Andrew G. Barto. Reinforcement learning: An introduction. MIT press, 2018.