

# Inverse Reinforcement Learning

Marius B. Mahiout and Herman V. Thieme

May 2023

# Introduction



How can we make an agent do a task?

---

<sup>1</sup>Stephen Adams, Tyler Cody, and Peter A Beling. “A survey of inverse reinforcement learning”. In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# Introduction



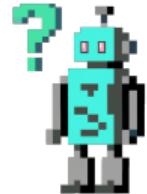
How can we make an agent do a task?

- Hard coded policy

---

<sup>1</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# Introduction



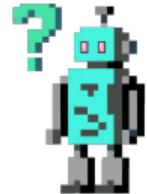
How can we make an agent do a task?

- Hard coded policy
  - Must anticipate correct actions in all future states<sup>1</sup>

---

<sup>1</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# Introduction



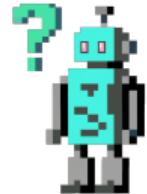
How can we make an agent do a task?

- Hard coded policy
  - Must anticipate correct actions in all future states<sup>1</sup>
- Reinforcement learning (RL)

---

<sup>1</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# Introduction



How can we make an agent do a task?

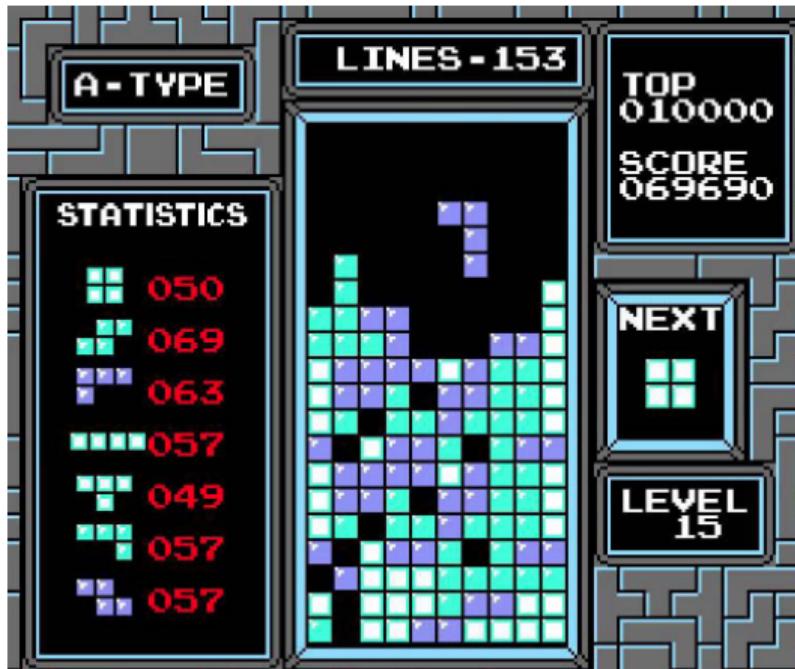
- Hard coded policy
  - Must anticipate correct actions in all future states<sup>1</sup>
- Reinforcement learning (RL)
  - RL requires a **reward function**

---

<sup>1</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

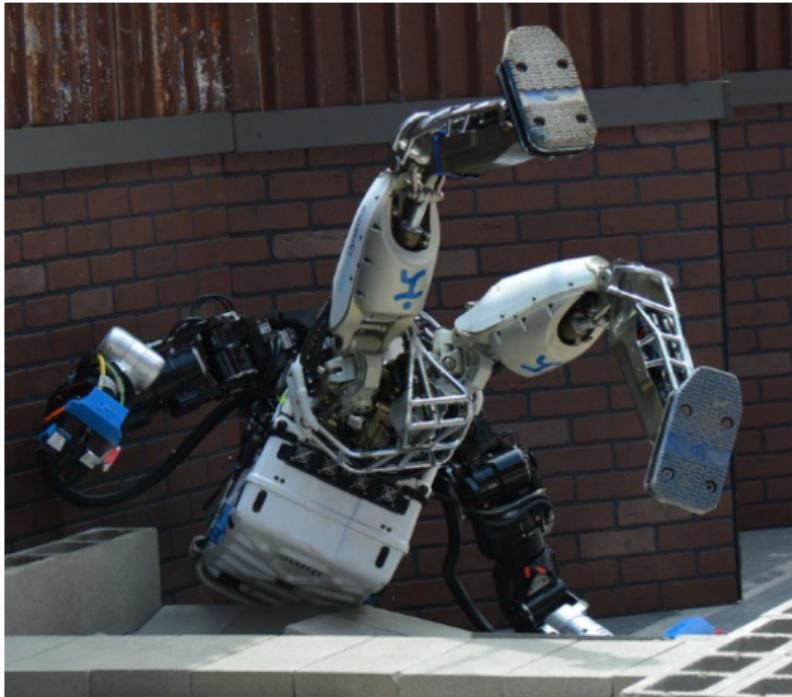
# The reward function

For some tasks, rewards are obvious

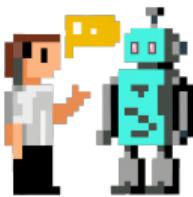


# The reward function

In other applications, not so obvious



# An alternative approach

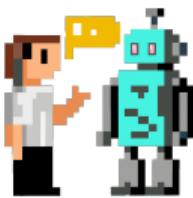


- Apprenticeship learning

---

<sup>2</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# An alternative approach

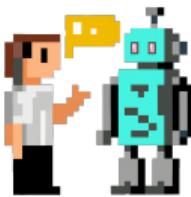


- Apprenticeship learning
- Direct policy learning

---

<sup>2</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# An alternative approach



- Apprenticeship learning
- Direct policy learning
  - Don't have rewards → can't do RL
  - Do supervised classification instead<sup>2</sup>

---

<sup>2</sup>Stephen Adams, Tyler Cody, and Peter A Beling. "A survey of inverse reinforcement learning". In: *Artificial Intelligence Review* 55.6 (2022), pp. 4307–4346.

# Rewards revisited

- Rewards are more generalizeable than policies

---

<sup>3</sup> Josep Call and Michael Tomasello. "Does the chimpanzee have a theory of mind? 30 years later". In: *Trends in cognitive sciences* 12.5 (2008), pp. 187–192.

# Rewards revisited

- Rewards are more generalizeable than policies
- Humans (and other animals<sup>3</sup>) can infer intentions!

---

<sup>3</sup> Josep Call and Michael Tomasello. "Does the chimpanzee have a theory of mind? 30 years later". In: *Trends in cognitive sciences* 12.5 (2008), pp. 187–192.

# Learning reward functions

- Instead of learning the policy, learn the reward function

---

<sup>4</sup>Saurabh Arora and Prashant Doshi. “A survey of inverse reinforcement learning: Challenges, methods and progress”. In: *Artificial Intelligence* 297 (2021), p. 103500.

# Learning reward functions

- Instead of learning the policy, learn the reward function
- This is the problem of inverse reinforcement learning (IRL)<sup>4</sup>

---

<sup>4</sup>Saurabh Arora and Prashant Doshi. "A survey of inverse reinforcement learning: Challenges, methods and progress". In: *Artificial Intelligence* 297 (2021), p. 103500.

# Learning reward functions

- Instead of learning the policy, learn the reward function
- This is the problem of inverse reinforcement learning (IRL)<sup>4</sup>
  - learned rewards → policy

---

<sup>4</sup>Saurabh Arora and Prashant Doshi. "A survey of inverse reinforcement learning: Challenges, methods and progress". In: *Artificial Intelligence* 297 (2021), p. 103500.

# Learning reward functions

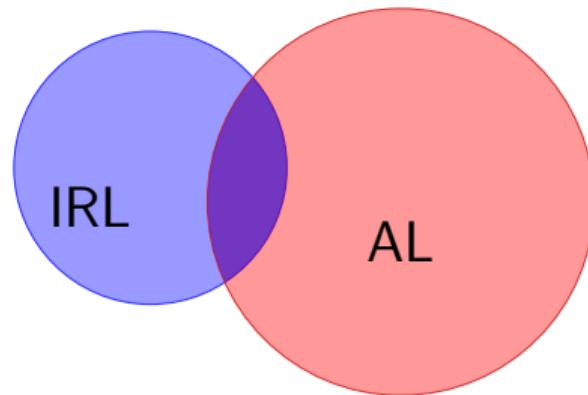
- Instead of learning the policy, learn the reward function
- This is the problem of inverse reinforcement learning (IRL)<sup>4</sup>
  - learned rewards → policy
  - $AL \cap IRL \neq \emptyset$

---

<sup>4</sup>Saurabh Arora and Prashant Doshi. "A survey of inverse reinforcement learning: Challenges, methods and progress". In: *Artificial Intelligence* 297 (2021), p. 103500.

# Learning reward functions

- Instead of learning the policy, learn the reward function
- This is the problem of inverse reinforcement learning (IRL)<sup>4</sup>
  - learned rewards → policy
  - $AL \cap IRL \neq \emptyset$
  - $IRL \not\subset AL$



---

<sup>4</sup>Saurabh Arora and Prashant Doshi. "A survey of inverse reinforcement learning: Challenges, methods and progress". In: *Artificial Intelligence* 297 (2021), p. 103500.

# This project

- Our project
  - Gridworld
  - DP
  - LP

# This project

- Our project
  - Gridworld
  - DP
  - LP
- Ng & Russell, 2000

# Infer reward from demonstrations



|   |   |   |   |   |
|---|---|---|---|---|
| ? | ? | ? | ? | ? |
| → | → | → | → | ? |
| ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? |
| ? | ? | ? | ? | ? |

# Infer reward from policies (Q)

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → | → |
| → | → | → | → | ↑ |
| → | → | → | ↑ | ↑ |
| ↑ | → | ↑ | ↑ | ↑ |
| → | → | ↑ | ↑ | ↑ |

**Exercise:** what is the right specification of the rewards in each state that produced the given policy

# Infer reward from policies (A)

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → | → |
| → | → | → | → | ↑ |
| → | → | → | ↑ | ↑ |
| ↑ | → | ↑ | ↑ | ↑ |
| → | → | ↑ | ↑ | ↑ |

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → | → |
| → | → | → | → | ↑ |
| → | → | → | ↑ | ↑ |
| ↑ | → | ↑ | ↑ | ↑ |
| → | → | ↑ | ↑ | ↑ |

# Infer reward from policies (A)

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → | → |
| → | → | → | → | ↑ |
| → | → | → | ↑ | ↑ |
| ↑ | → | ↑ | ↑ | ↑ |
| → | → | ↑ | ↑ | ↑ |

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → | → |
| → | → | → | → | ↑ |
| → | → | → | ↑ | ↑ |
| ↑ | → | ↑ | ↑ | ↑ |
| → | → | ↑ | ↑ | ↑ |

# Methods

## The forward and inverse problems<sup>5</sup>

### Forward problem

- Given:
  - states  $s \in \mathcal{S}$ , actions  $a \in \mathcal{A}$
  - transition probabilities  
 $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}, P_{s,a}(s')$
  - reward function  $R(s)$
- Goal:
  - learn optimal policy  $\pi(s)$

### Inverse problem

- Given:
  - states  $s \in \mathcal{S}$ , actions  $a \in \mathcal{A}$
  - transition probabilities  
 $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}, P_{s,a}(s')$
  - optimal policy  $\pi(s)$
- Goal:
  - learn reward function  $R(s)$

---

<sup>5</sup> Andrew Y Ng, Stuart Russell, et al. “Algorithms for inverse reinforcement learning.”. In: *Icml*. Vol. 1. 2000, p. 2.

# Methods

## The forward problem

- Value iteration<sup>6</sup>

---

<sup>6</sup>Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*.

# Methods

## The forward problem

- Value iteration<sup>6</sup>
- The Bellman equations:

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \quad (1)$$

$$Q^\pi(s, a) = R(s) + \gamma \sum_{s'} P_{s,a}(s') V^\pi(s') \quad (2)$$

---

<sup>6</sup>Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*.

# Methods

## The forward problem

- Value iteration<sup>6</sup>
- The Bellman equations:

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \quad (1)$$

$$Q^\pi(s, a) = R(s) + \gamma \sum_{s'} P_{s,a}(s') V^\pi(s') \quad (2)$$

- Optimality condition:

$$\forall s \in \mathcal{S}, \pi(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^\pi(s, a) \quad (3)$$

---

<sup>6</sup>Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*.

# Methods

## Characterization of reward functions

Which rewards  $\mathbf{R}$  are compatible with a given policy?

Answer:

$\mathbf{R}$  such that  $\forall s \in \mathcal{S}, a \in \mathcal{A} \setminus \pi(s)$ ,

$$(\mathbf{P}_{\pi(s)} - \mathbf{P}_a)(\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{R} \geq 0 \quad (4)$$

# Methods

## Proof of characterization

- Substituting (2) into (3) yields:  
 $\pi$  optimal  $\iff \forall s \in \mathcal{S}, \pi(s) \in \operatorname{argmax}_a \sum_{s'} P_{s,a}(s') V^\pi(s')$

# Methods

## Proof of characterization

- Substituting (2) into (3) yields:

$$\begin{aligned}\pi \text{ optimal} &\iff \forall s \in \mathcal{S}, \pi(s) \in \operatorname{argmax}_a \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, a \in \mathcal{A}, \\ &\quad \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \geq \sum_{s'} P_{s,a}(s') V^\pi(s')\end{aligned}$$

# Methods

## Proof of characterization

- Substituting (2) into (3) yields:

$$\begin{aligned}\pi \text{ optimal} &\iff \forall s \in \mathcal{S}, \pi(s) \in \operatorname{argmax}_a \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, a \in \mathcal{A}, \\ &\quad \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \geq \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \setminus \pi(s), \mathbf{P}_{\pi(s)} \mathbf{V}^\pi \geq \mathbf{P}_a \mathbf{V}^\pi\end{aligned}$$

# Methods

## Proof of characterization

- Substituting (2) into (3) yields:

$$\begin{aligned}\pi \text{ optimal} &\iff \forall s \in \mathcal{S}, \pi(s) \in \operatorname{argmax}_a \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, a \in \mathcal{A}, \\ &\quad \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \geq \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \setminus \pi(s), \mathbf{P}_{\pi(s)} \mathbf{V}^\pi \geq \mathbf{P}_a \mathbf{V}^\pi\end{aligned}$$

- Rewriting (1) as:

$$\mathbf{V}^\pi = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{R},$$

# Methods

## Proof of characterization

- Substituting (2) into (3) yields:

$$\begin{aligned}\pi \text{ optimal} &\iff \forall s \in \mathcal{S}, \pi(s) \in \operatorname{argmax}_a \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, a \in \mathcal{A}, \\ &\quad \sum_{s'} P_{s,\pi(s)}(s') V^\pi(s') \geq \sum_{s'} P_{s,a}(s') V^\pi(s') \\ &\iff \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \setminus \pi(s), \mathbf{P}_{\pi(s)} \mathbf{V}^\pi \geq \mathbf{P}_a \mathbf{V}^\pi\end{aligned}$$

- Rewriting (1) as:

$$\mathbf{V}^\pi = (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{R},$$

- we obtain:

$$\mathbf{P}_{\pi(s)} (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{R} \geq \mathbf{P}_a (\mathbf{I} - \gamma \mathbf{P}_\pi)^{-1} \mathbf{R}.$$

# Methods

## The maximum margin principle

- Multiple  $\mathbf{R}$ 's that satisfy the condition (4)

---

<sup>7</sup> Andrew Y Ng, Stuart Russell, et al. "Algorithms for inverse reinforcement learning.". In: *Icml*. Vol. 1. 2000, p. 2.

# Methods

## The maximum margin principle

- Multiple  $\mathbf{R}$ 's that satisfy the condition (4)
- Heuristic: make the expert's policy look "maximally better"  
than all other policies<sup>7</sup>

---

<sup>7</sup> Andrew Y Ng, Stuart Russell, et al. "Algorithms for inverse reinforcement learning.". In: *Icml*. Vol. 1. 2000, p. 2.

# Methods

## The maximum margin principle

- Multiple  $\mathbf{R}$ 's that satisfy the condition (4)
- Heuristic: make the expert's policy look "maximally better" than all other policies<sup>7</sup>
- Maximize the *margin* between the values of expert policy and the next best actions:

$$\sum_{s \in \mathcal{S}} \left[ Q^\pi(s, \pi(s)) - \max_{a \neq \pi(s)} Q^\pi(s, a) \right]$$

---

<sup>7</sup> Andrew Y Ng, Stuart Russell, et al. "Algorithms for inverse reinforcement learning." . In: *Icml*. Vol. 1. 2000, p. 2.

# Methods

## The inverse problem

$$\max_{\mathbf{R}} \left[ \sum_{s=1}^N \min_{a \neq \pi(s)} \{ (\mathbf{P}_{\pi(s)}(s) - \mathbf{P}_a(s))(\mathbf{I} - \mathbf{P}_{\pi(s)})^{-1} \mathbf{R} \} - \lambda \|\mathbf{R}\|_1 \right]$$
$$s.t. \quad \forall s \in \mathcal{S}, a \neq \pi(s), \quad (\mathbf{P}_{\pi(s)} - \mathbf{P}_a)(\mathbf{I} - \gamma \mathbf{P}_{\pi})^{-1} \mathbf{R} \geq 0$$

# Methods

## The inverse problem

$$\max_{\mathbf{R}} \left[ \sum_{s=1}^N \min_{a \neq \pi(s)} \{ (\mathbf{P}_{\pi(s)}(s) - \mathbf{P}_a(s))(\mathbf{I} - \mathbf{P}_{\pi(s)})^{-1} \mathbf{R} \} - \lambda \|\mathbf{R}\|_1 \right]$$
$$s.t. \quad \forall s \in \mathcal{S}, a \neq \pi(s), \quad (\mathbf{P}_{\pi(s)} - \mathbf{P}_a)(\mathbf{I} - \gamma \mathbf{P}_{\pi})^{-1} \mathbf{R} \geq 0$$

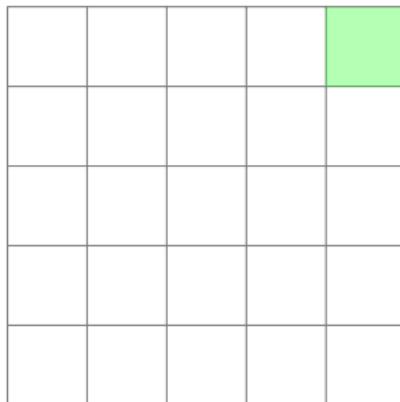
# Methods

## The inverse problem

$$\max_{\mathbf{R}} \left[ \sum_{s=1}^N \min_{a \neq \pi(s)} \{ (\mathbf{P}_{\pi(s)}(s) - \mathbf{P}_a(s))(\mathbf{I} - \mathbf{P}_{\pi(s)})^{-1} \mathbf{R} \} - \lambda \|\mathbf{R}\|_1 \right.$$
$$\left. s.t. \quad \forall s \in \mathcal{S}, a \neq \pi(s), \quad (\mathbf{P}_{\pi(s)} - \mathbf{P}_a)(\mathbf{I} - \gamma \mathbf{P}_{\pi})^{-1} \mathbf{R} \geq 0 \right]$$

# Experiment

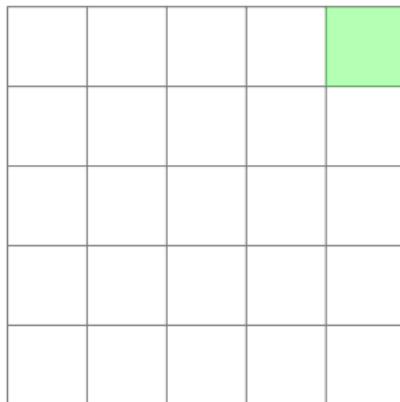
Ng and Russell



- Green: Absorbing state

# Experiment

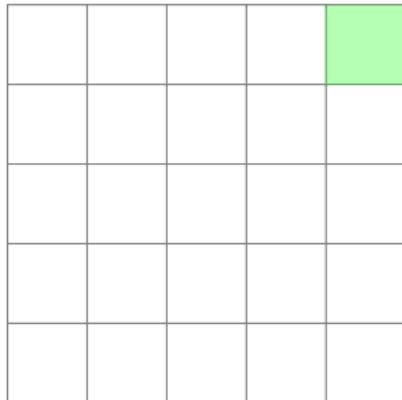
Ng and Russell



- Green: Absorbing state
- $R((1, 5)) = 1$  else  $R(s) = 0$

# Experiment

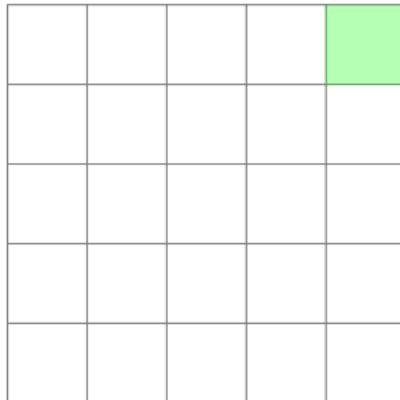
Ng and Russell



- Green: Absorbing state
- $R((1, 5)) = 1$  else  $R(s) = 0$
- Actions =  
*North, South, East, West*

# Experiment

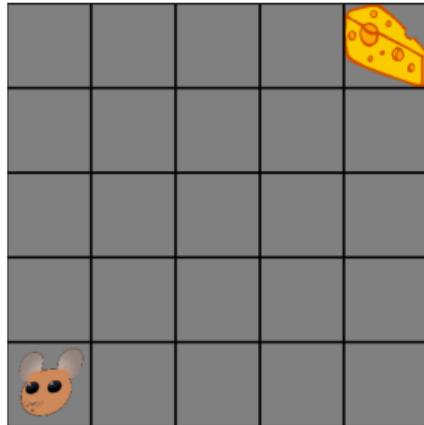
Ng and Russell



- Green: Absorbing state
- $R((1, 5)) = 1$  else  $R(s) = 0$
- Actions =  
*North, South, East, West*
- Noise = 0.3

# Experiment

Replication<sup>8</sup>



---

<sup>8</sup>Greg Brockman et al. *OpenAI Gym*. 2016. eprint: arXiv:1606.01540.

# Experiment

## Policies

Ng and Russell

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → |   |
| ↑ | → | → | ↑ | ↑ |
| ↑ | ↑ | ↑ | ↑ | ↑ |
| ↑ | ↑ | → | ↑ | ↑ |
| ↑ | → | → | → | ↑ |

# Experiment

## Policies

Ng and Russell

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → |   |
| ↑ | → | → | ↑ | ↑ |
| ↑ | ↑ | ↑ | ↑ | ↑ |
| ↑ | ↑ | → | ↑ | ↑ |
| ↑ | → | → | → | ↑ |

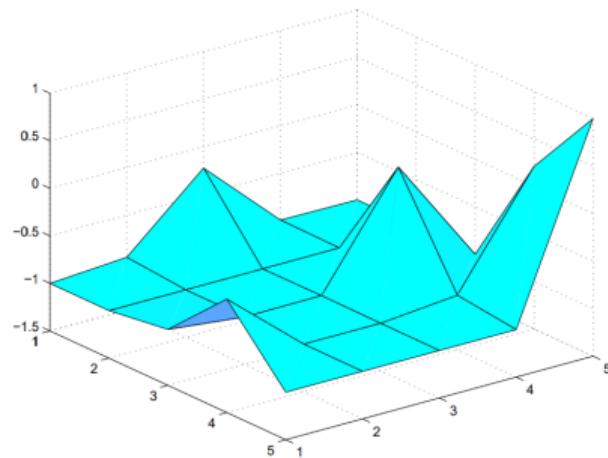
Replication

|   |   |   |   |   |
|---|---|---|---|---|
| → | → | → | → |   |
| → | → | → | → | ↑ |
| → | → | → | ↑ | ↑ |
| ↑ | → | ↑ | ↑ | ↑ |
| → | → | ↑ | ↑ | ↑ |

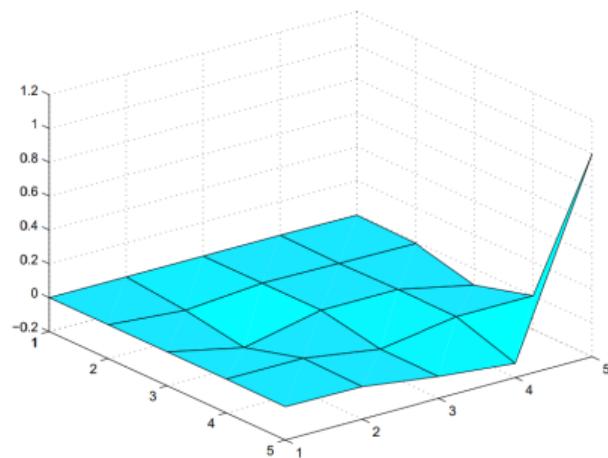
# Experiment

Ng and Russell

$$\lambda = 0$$



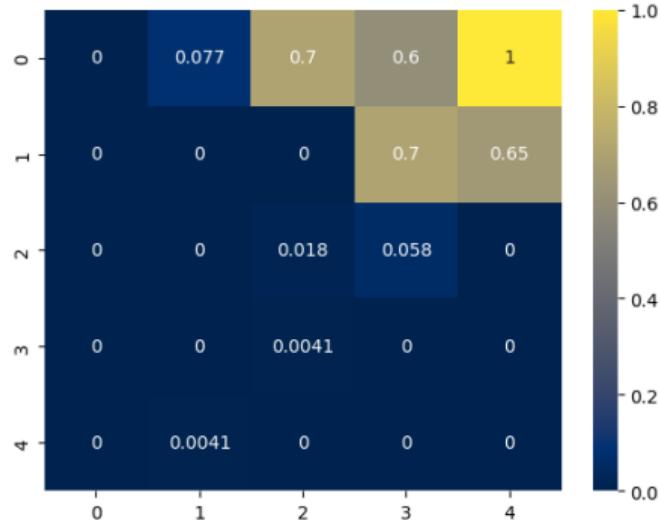
$$\lambda = 1.05$$



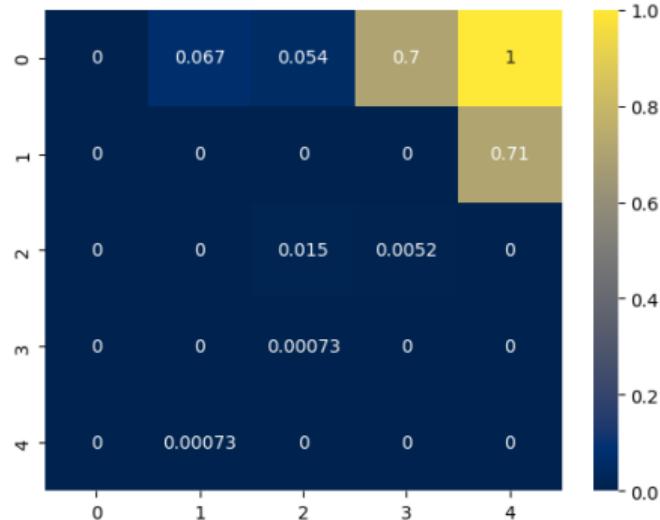
# Experiment

## Replication

$\lambda = 0$

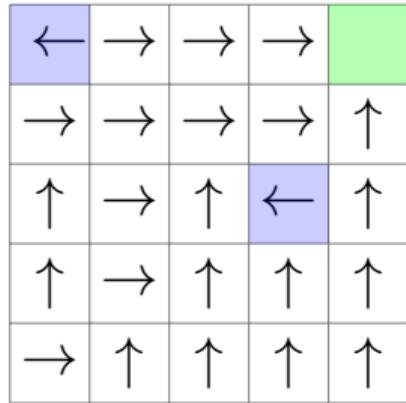


$\lambda = 1.05$



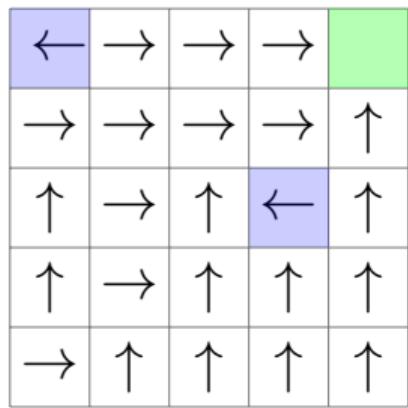
# Experiment

Variation: Suboptimal policy

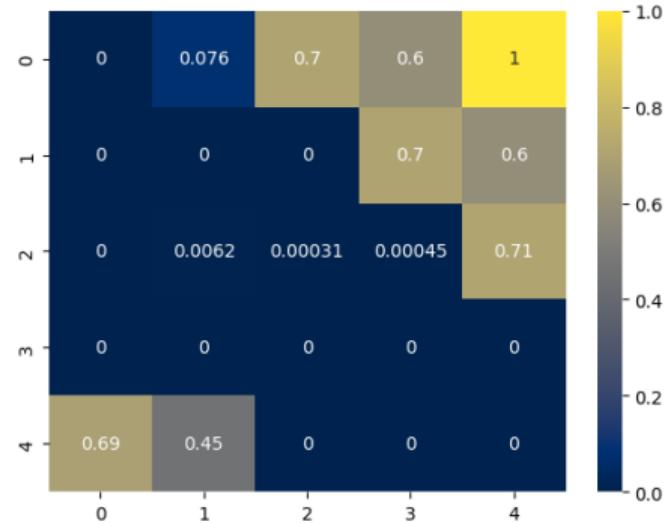


# Experiment

Variation: Suboptimal policy

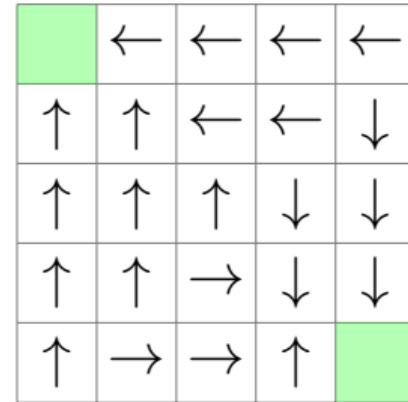
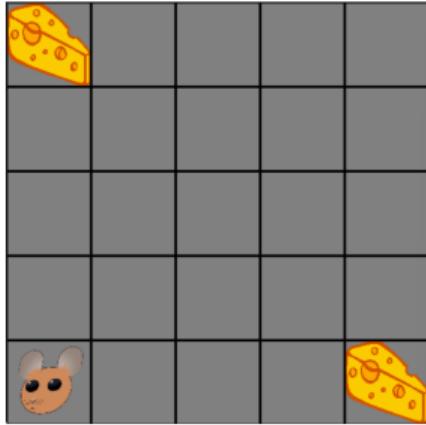


$$\lambda = 1.05, \epsilon = 0.2$$



# Experiment

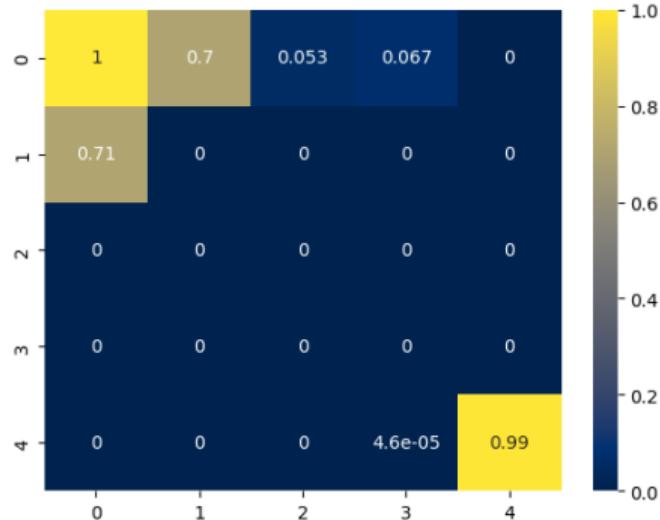
Variation: Two absorbing states



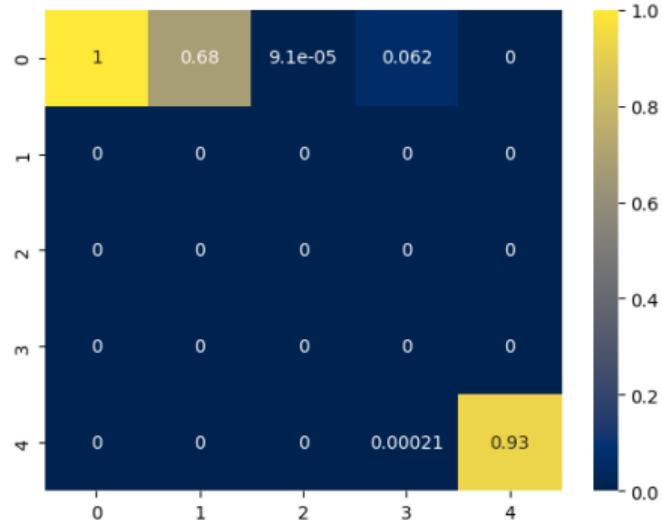
# Experiment

Variation: Two absorbing states

$\lambda = 0$



$\lambda = 1.05$



# Conclusion

## Summary

- Refining objectives

# Conclusion

## Summary

- Refining objectives
- Reward ambiguity

# Conclusion

## Summary

- Refining objectives
- Reward ambiguity
- Solution strategy

$$\max_{\mathbf{R}} \left[ \sum_{s=1}^N \min_{a \neq \pi(s)} \{ (\mathbf{P}_{\pi(s)}(s) - \mathbf{P}_a(s))(\mathbf{I} - \mathbf{P}_{\pi(s)})^{-1} \mathbf{R} \} - \lambda \|\mathbf{R}\|_1 \right]$$
$$s.t. \quad \forall s \in \mathcal{S}, a \neq \pi(s), \quad (\mathbf{P}_{\pi(s)} - \mathbf{P}_a)(\mathbf{I} - \gamma \mathbf{P}_{\pi})^{-1} \mathbf{R} \geq 0$$

# Conclusion

## Future directions

- Limitations
  - Finite state spaces
  - Known transition probabilities
  - Known optimal policy

# Conclusion

## Future directions

- Limitations
  - Finite state spaces
  - Known transition probabilities
  - Known optimal policy
- Other methods
  - LP linear approximation for continuous state space
  - LP from sampled trajectories
  - Maximum entropy methods