

Course 3, Module 1

On-policy Prediction with

Approximation

CMPUT 365
Fall 2021

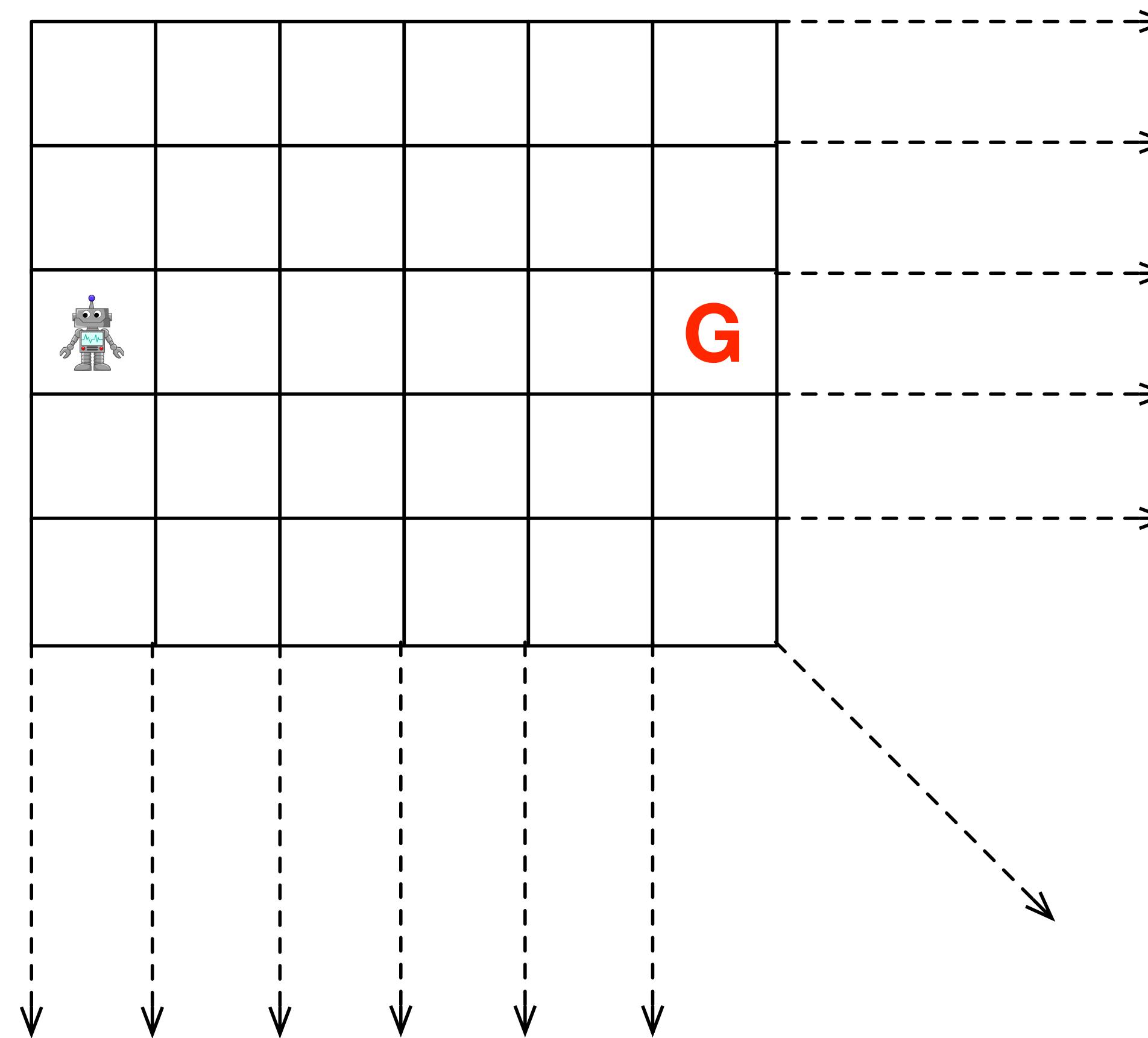
Announcements

- **Reading week next week**
 - *No office hours with Adam. I am taking vacation*
- **Midterm when you come back from reading week**
 - practice midterm on eClass

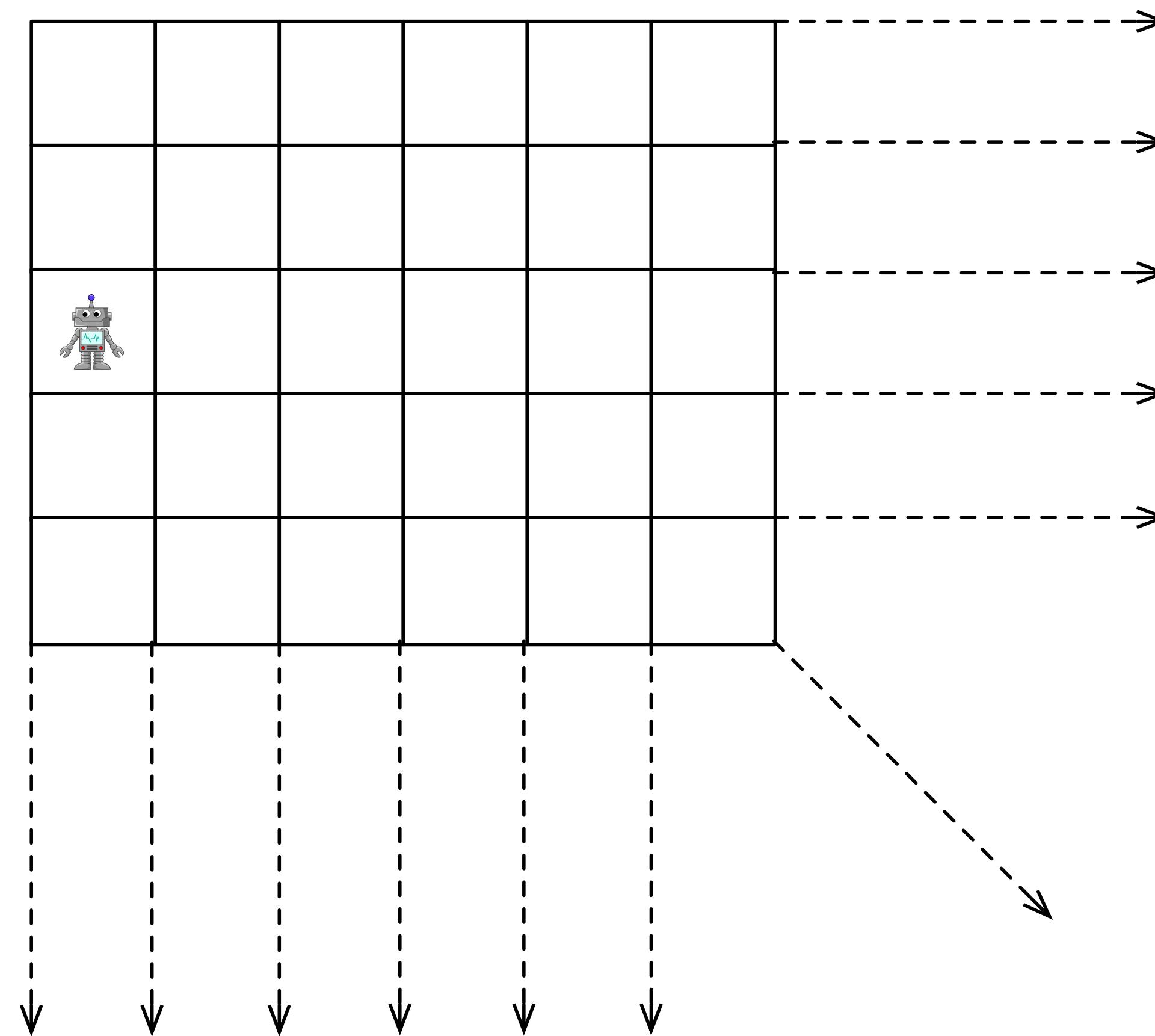
Moving to Approximation

- Our goal remains the same, as in Course 1 and Course 2
- But now we cannot represent value functions perfectly
 - because the space is too big
- Course 3 is about how to extend our algorithms to approximate value functions

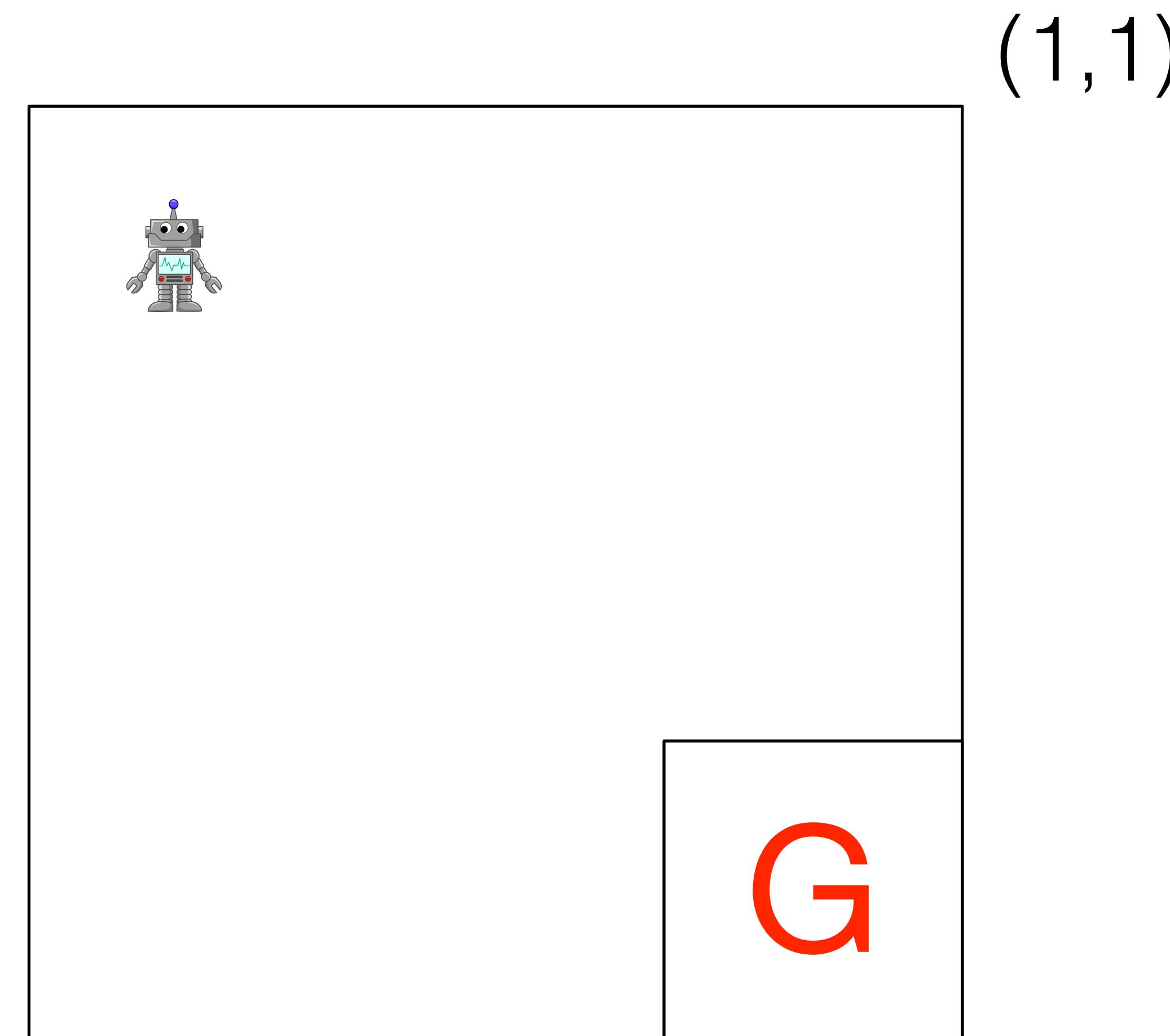
Imagine a huge state space



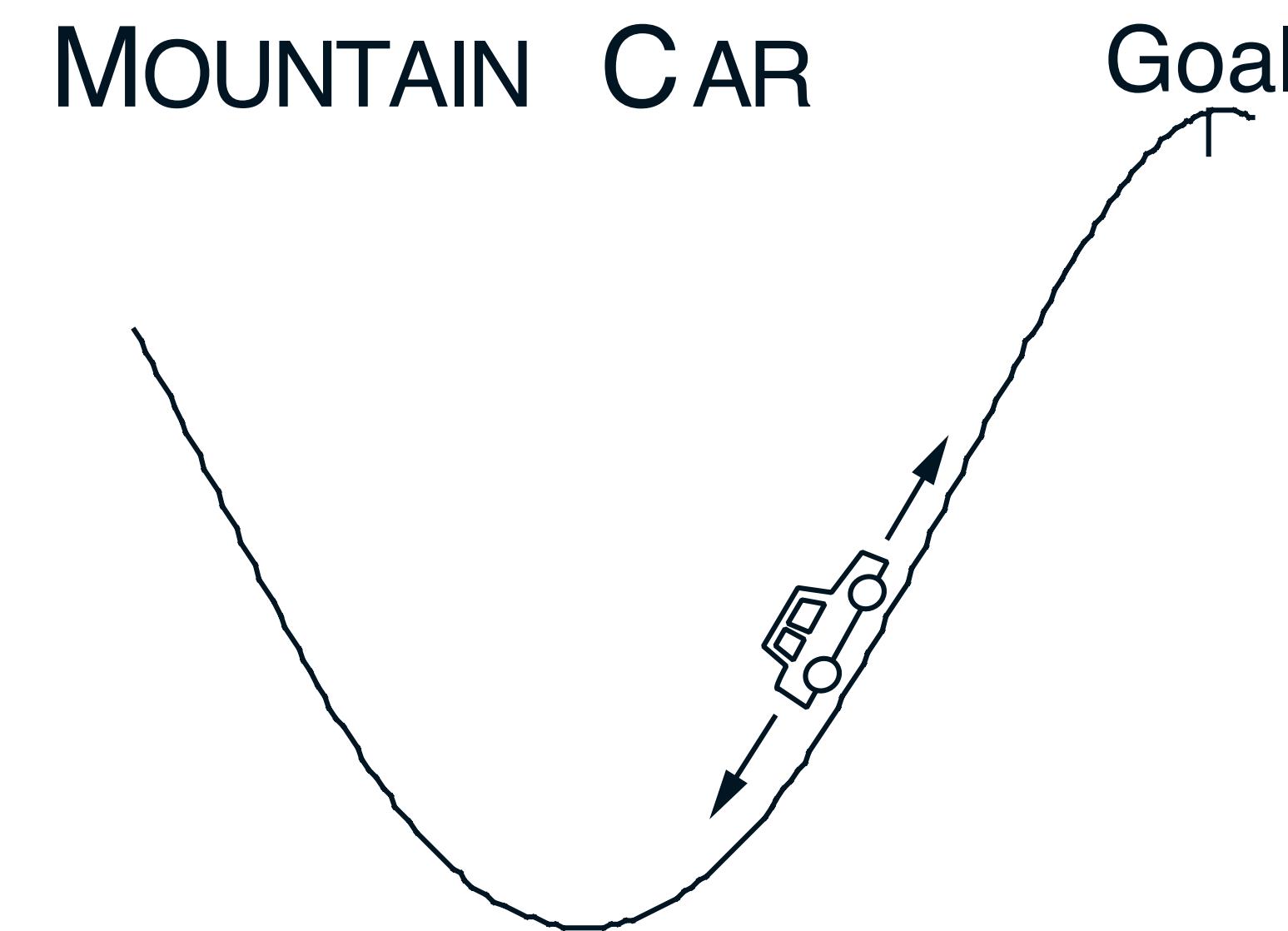
Imagine a huge state space



Imagine a continuous state space



Another continuous state domain



$$p_{t+1} \doteq \text{bound}[p_t + \dot{p}_{t+1}]$$

$$\dot{p}_{t+1} \doteq \text{bound}[\dot{p}_t + 0.001A_t - 0.0025 \cos(3p_t)]$$

Video 1: Moving to Parameterized Functions

- **Using parameterized functions to represent value functions.** From tables of values to more general functions over states
- Goals:
 - Understand how we can use parameterized functions to approximate values.
 - Explain linear value function approximation.
 - Recognize that the tabular case is a special case of linear value function approximation
 - Understand that there are many ways to parameterize an approximate value function.

$$\cancel{V(s) \approx} v_{\pi}(s)$$

1.71

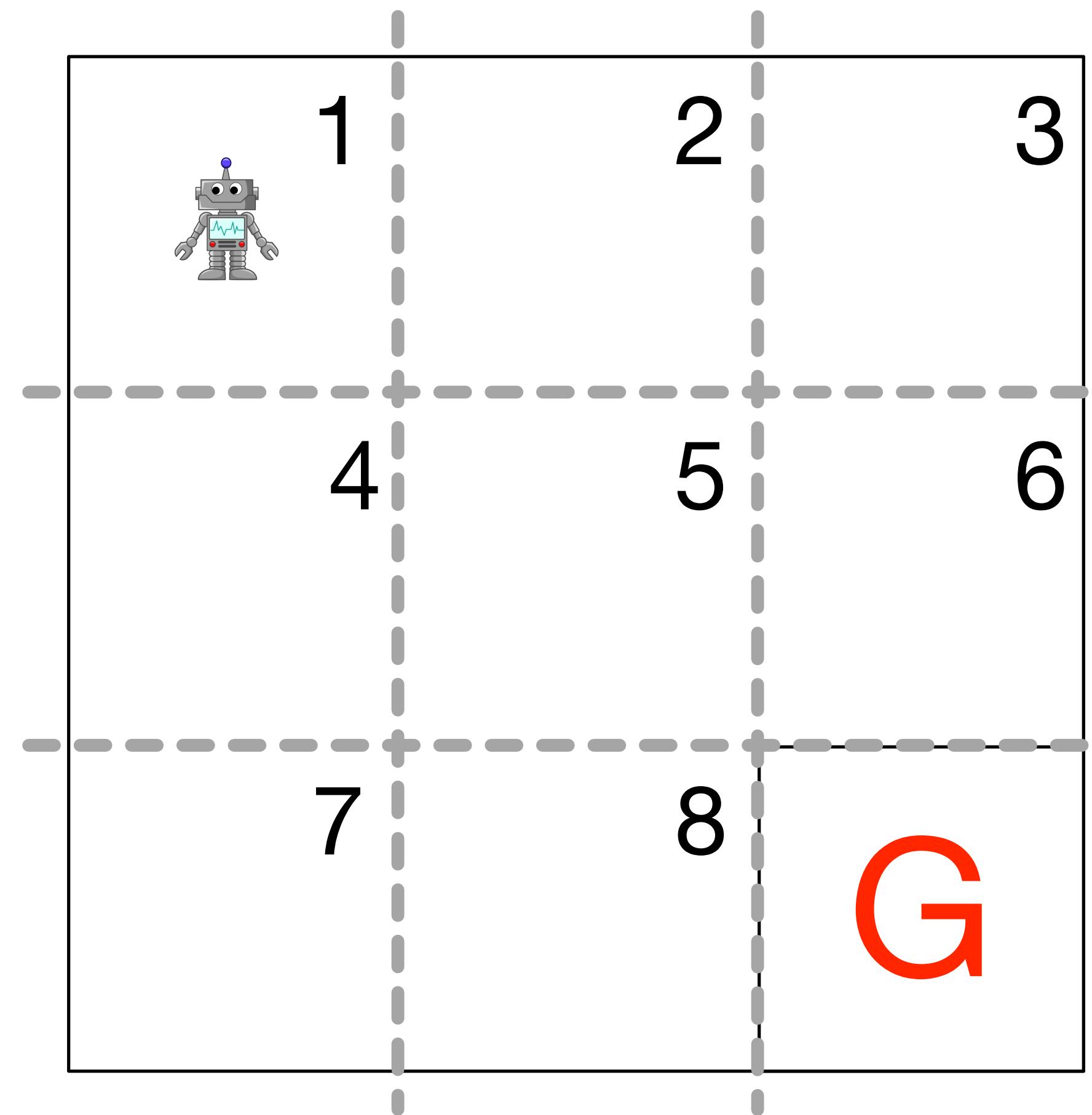
$$\mathbf{w} \in \mathbb{R}^d, \text{ e.g., } \mathbf{w} =$$

**parameter
vector**

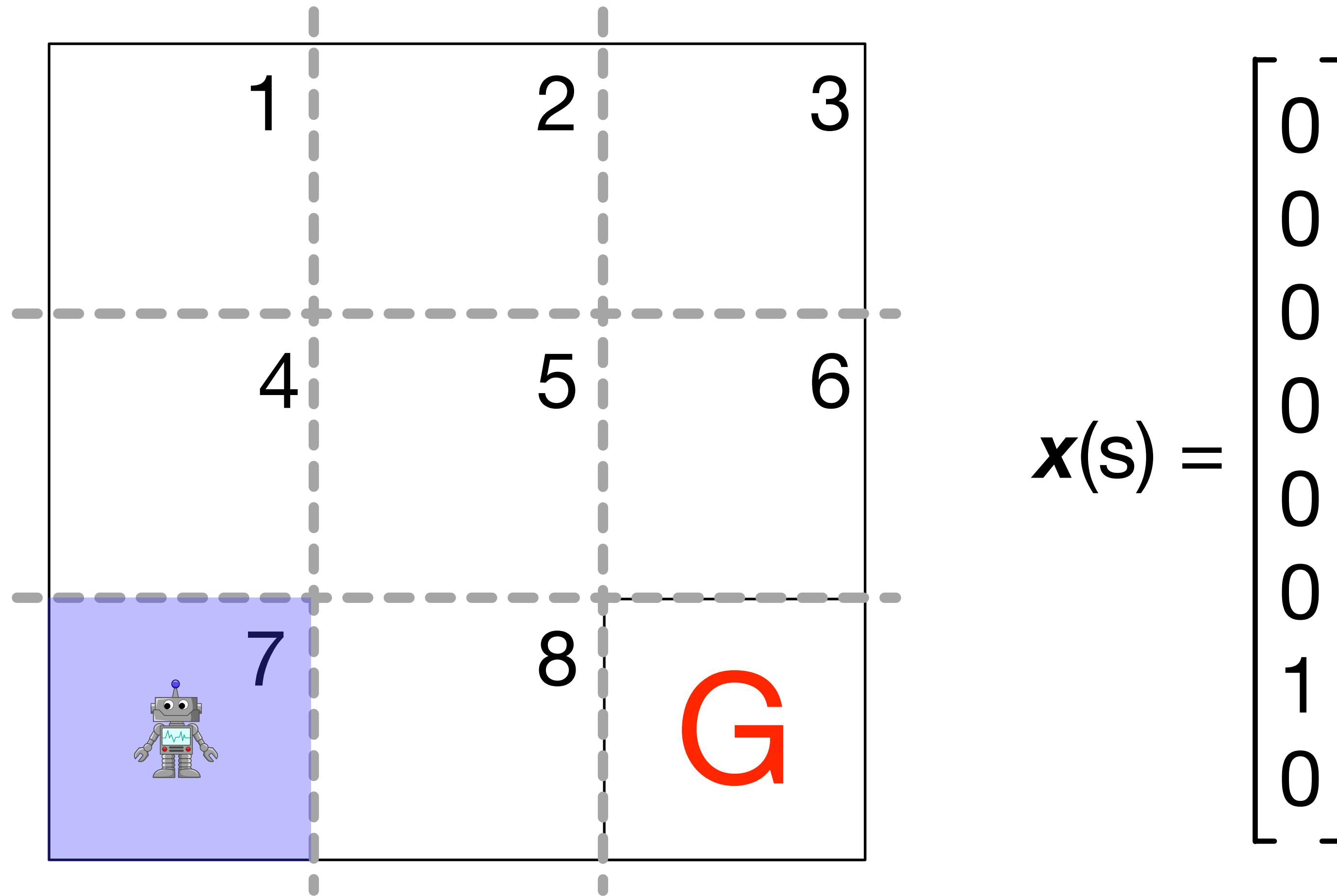
$$\mathbf{w} = \begin{bmatrix} 2.1 \\ 0.01 \\ -1.1 \\ 1.2 \\ -0.1 \\ 0.01 \\ 4.93 \\ 0.5 \end{bmatrix}, \quad \mathbf{x}(s) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \quad \mathbf{x} : \mathcal{S} \rightarrow \mathbb{R}^d$$

**feature
vector**

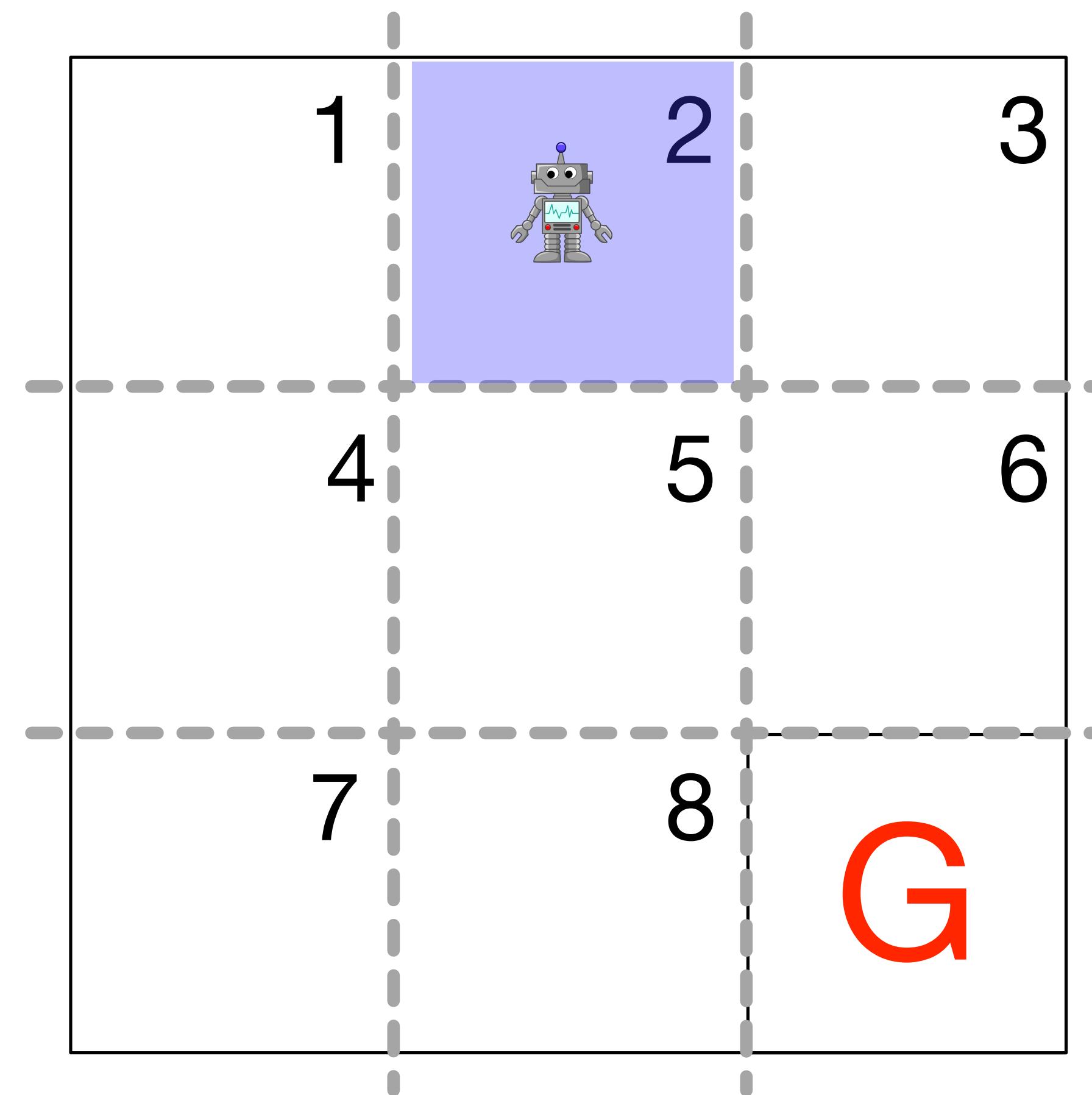
Let's look at a simple state aggregation



Here is the feature vector if the agent is somewhere in the bottom left

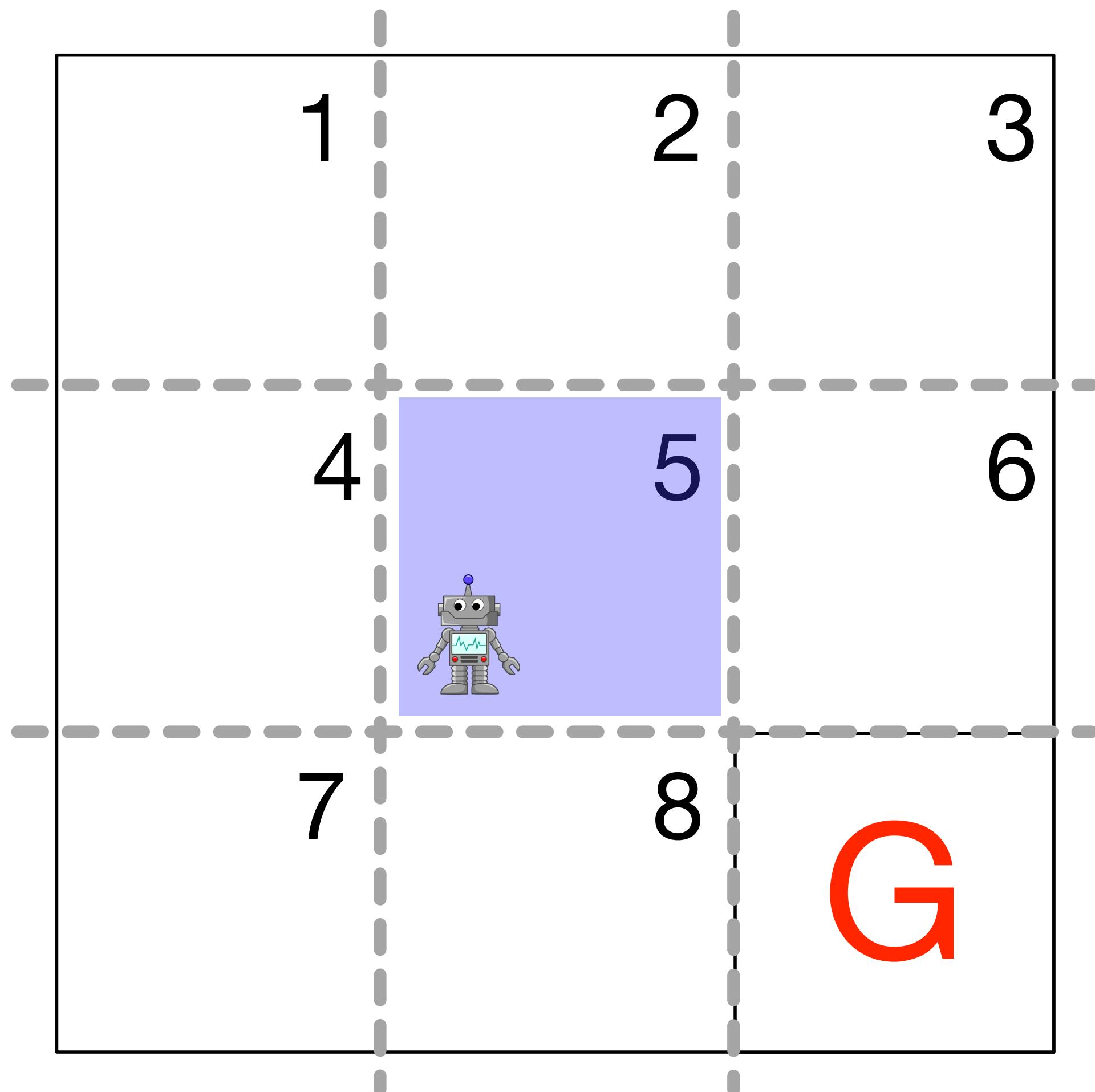


What would the feature vector be if the agent was somewhere in the top middle?



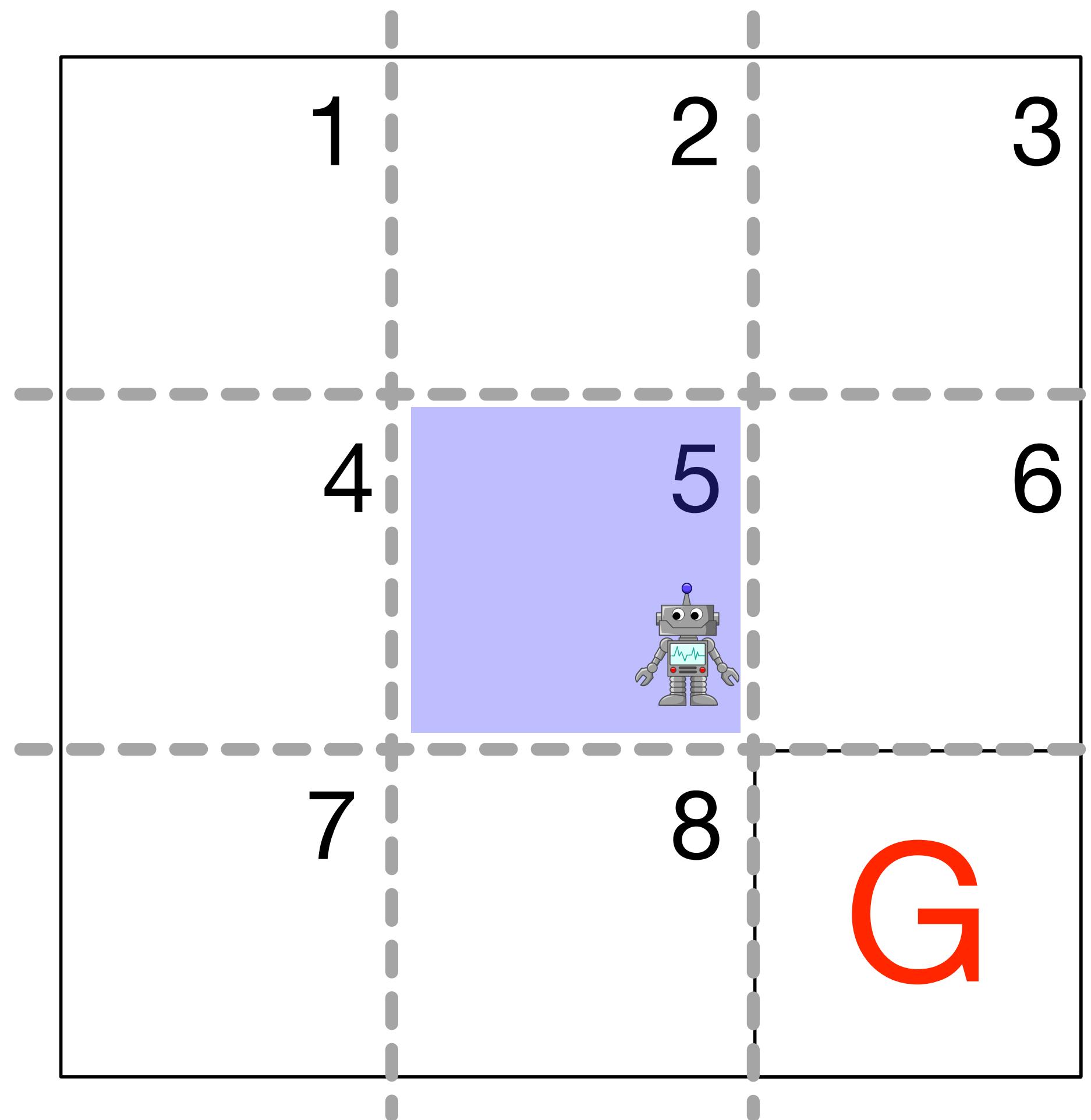
$$\mathbf{x}(s) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

What would the feature vector be if the agent was in the middle?



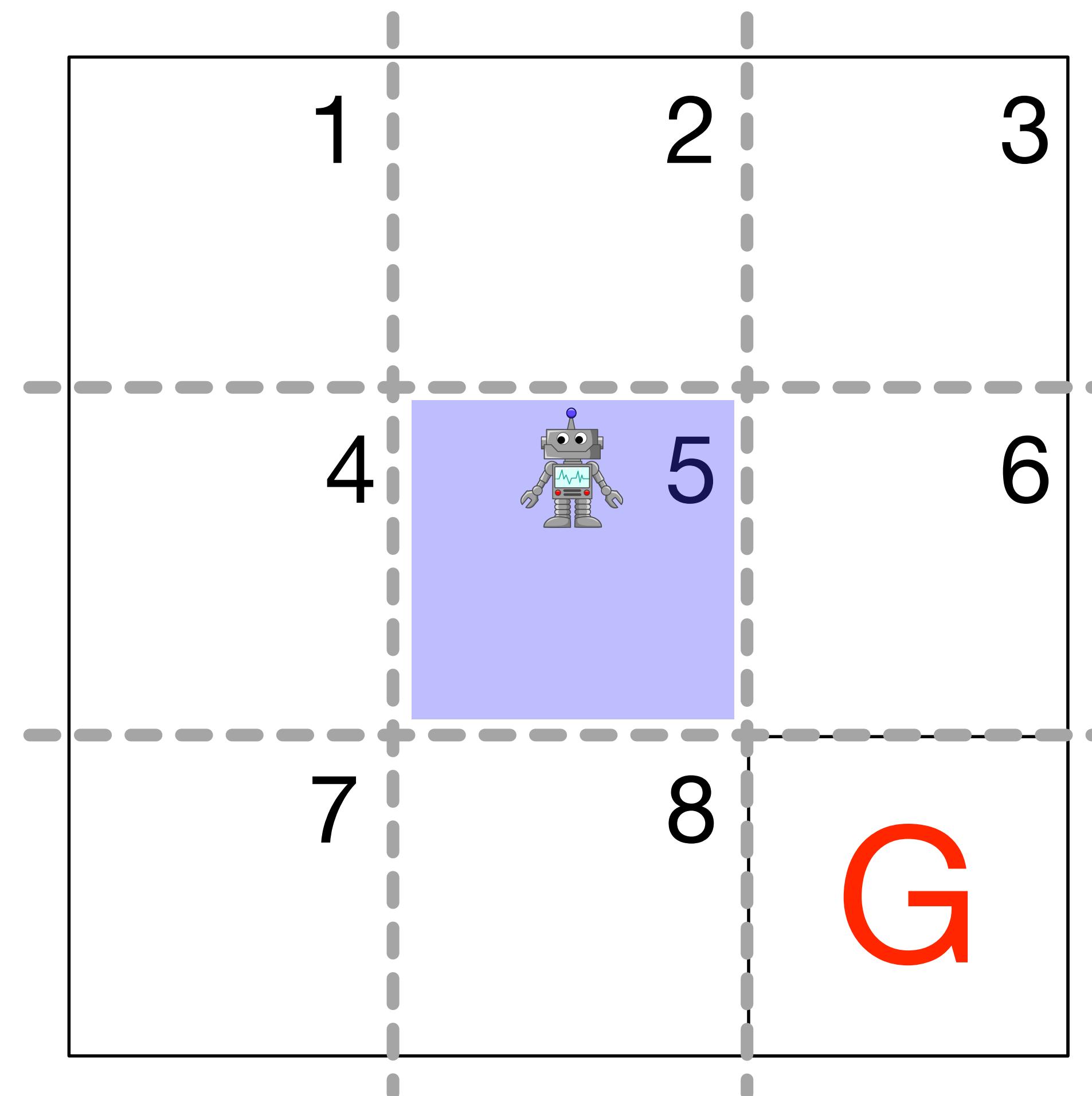
$$x(s) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

How about here?



$$x(s) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

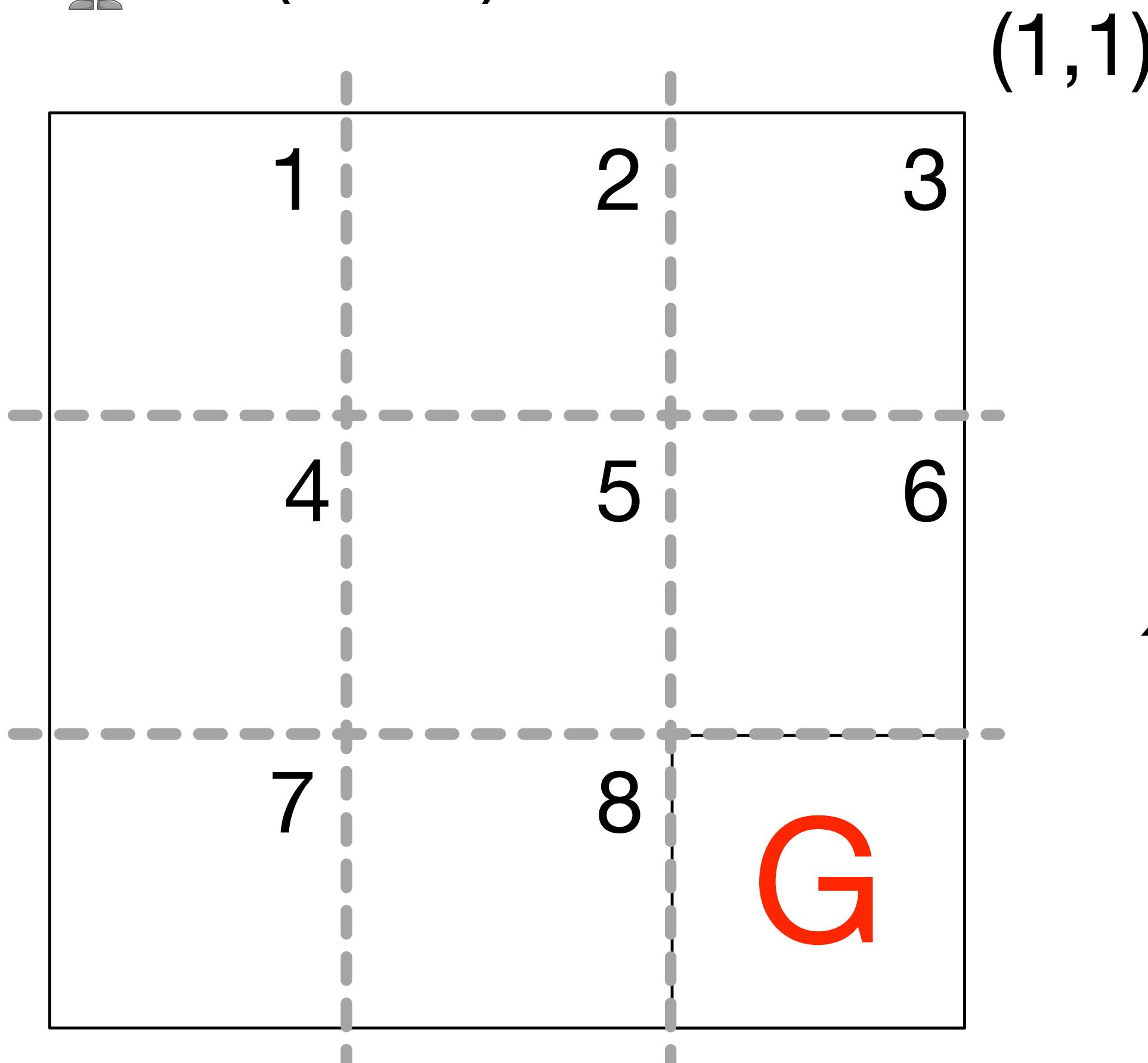
Or here?



$$x(s) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

What would the feature vector be if the agent was at this point?

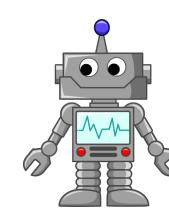
$$\text{Robot} = (.8, .8)$$

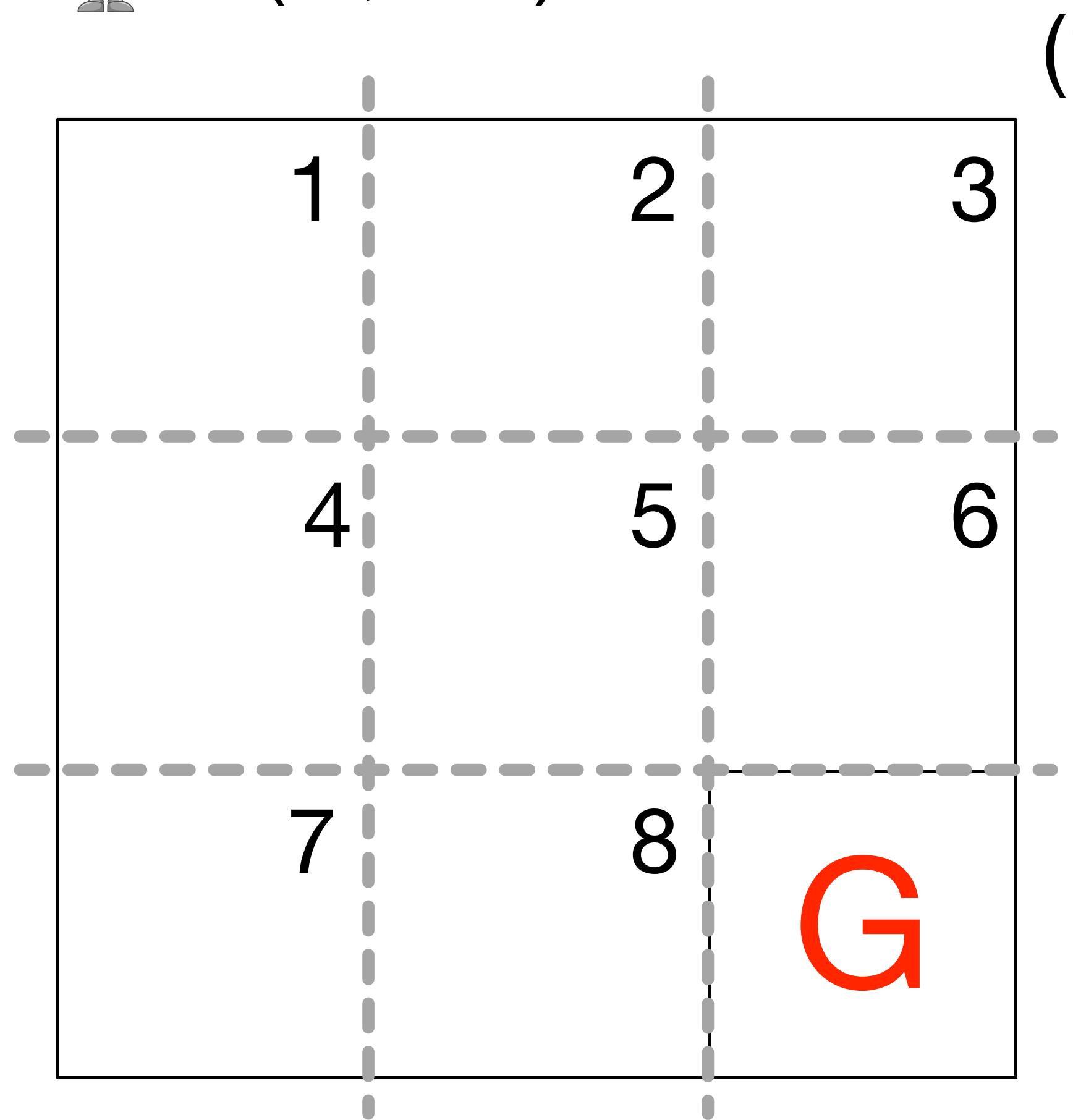


$$x(s) = \text{Robot} = x(0.8, 0.8) =$$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Or here?

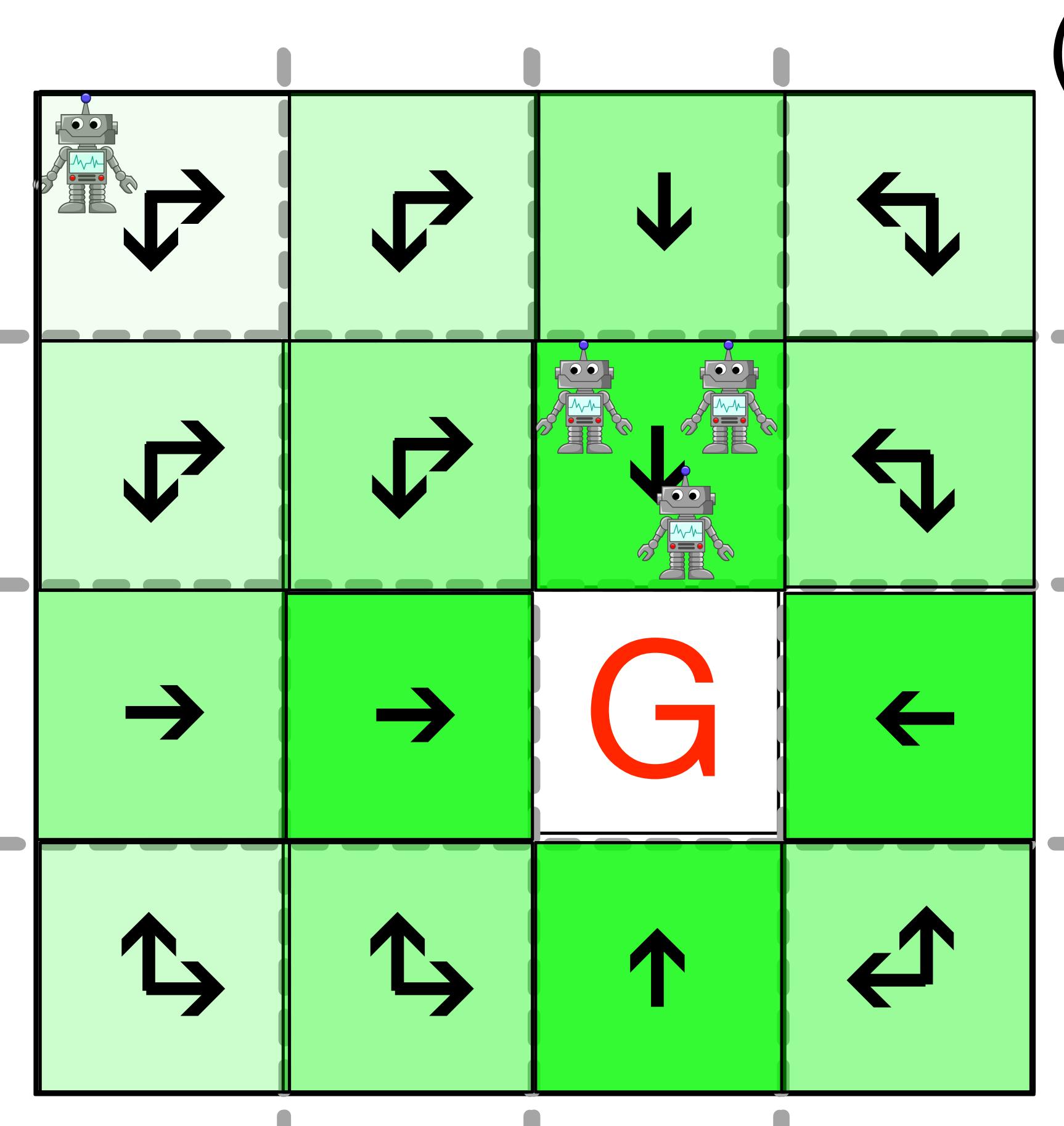
 = (.1, .45)



$x(s) = \begin{array}{c} \text{robot icon} \end{array} = x(0.1, 0.45) =$

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

What might the final estimate of the value function look like with this state aggregation?



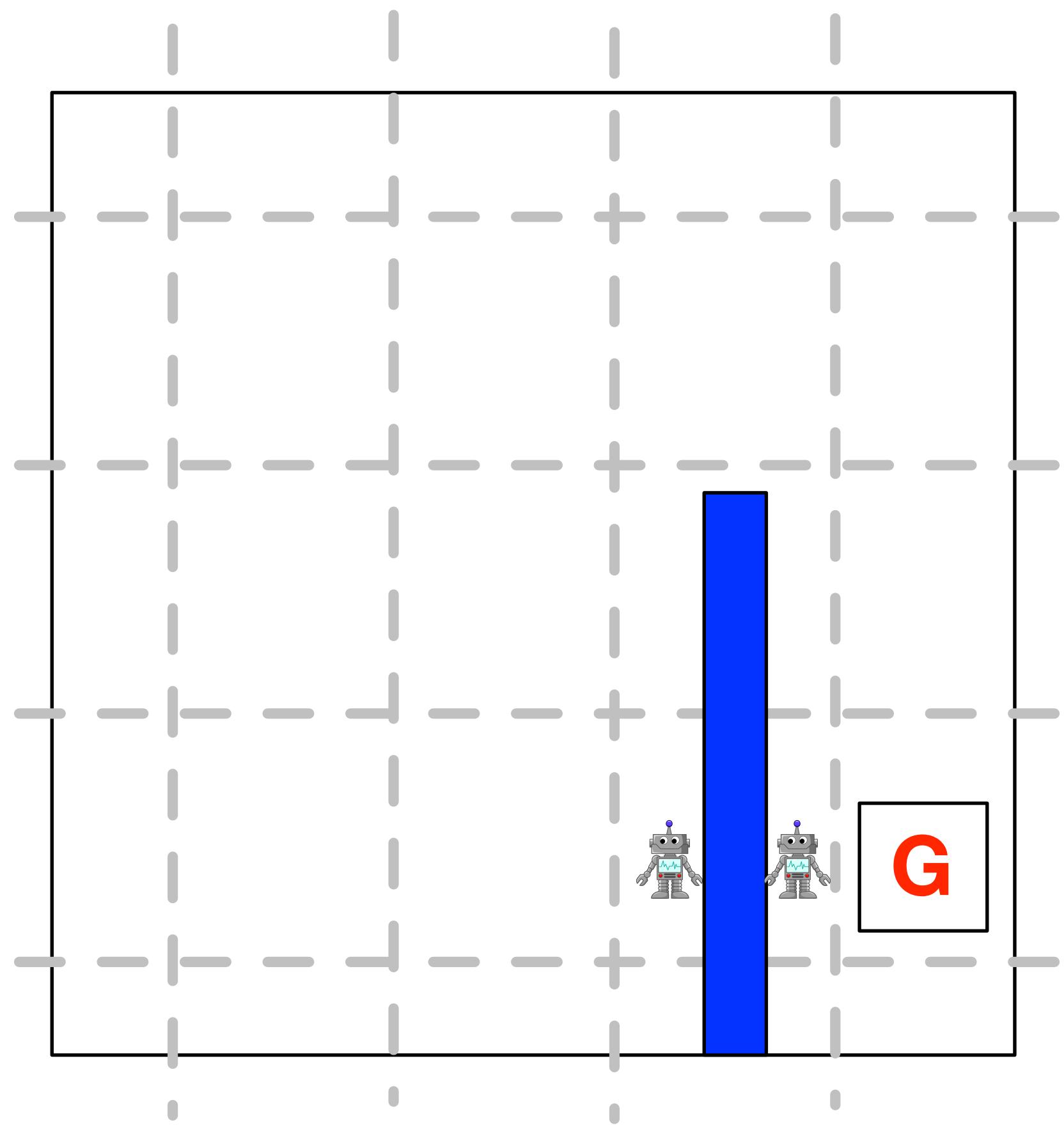
(1,1)

- $R = +1$ per step
- episodic, $\gamma = 1$
- agent starts in the top left corner
- π = shortest path policy
- what should $\hat{v}(s, w)$ look like?

Video 2: Generalization and Discrimination

- A key concept in machine learning. We cannot learn all the values separately (in fact we wouldn't want to), so we have to make choices.
- Goals:
 - Understand what is meant by generalization and discrimination
 - Understand how generalization can be beneficial
 - Explain why we want both generalization and discrimination from our function approximation

**Exercise: Is there any issue with this state aggregation?
Can we represent the optimal action-value function?**



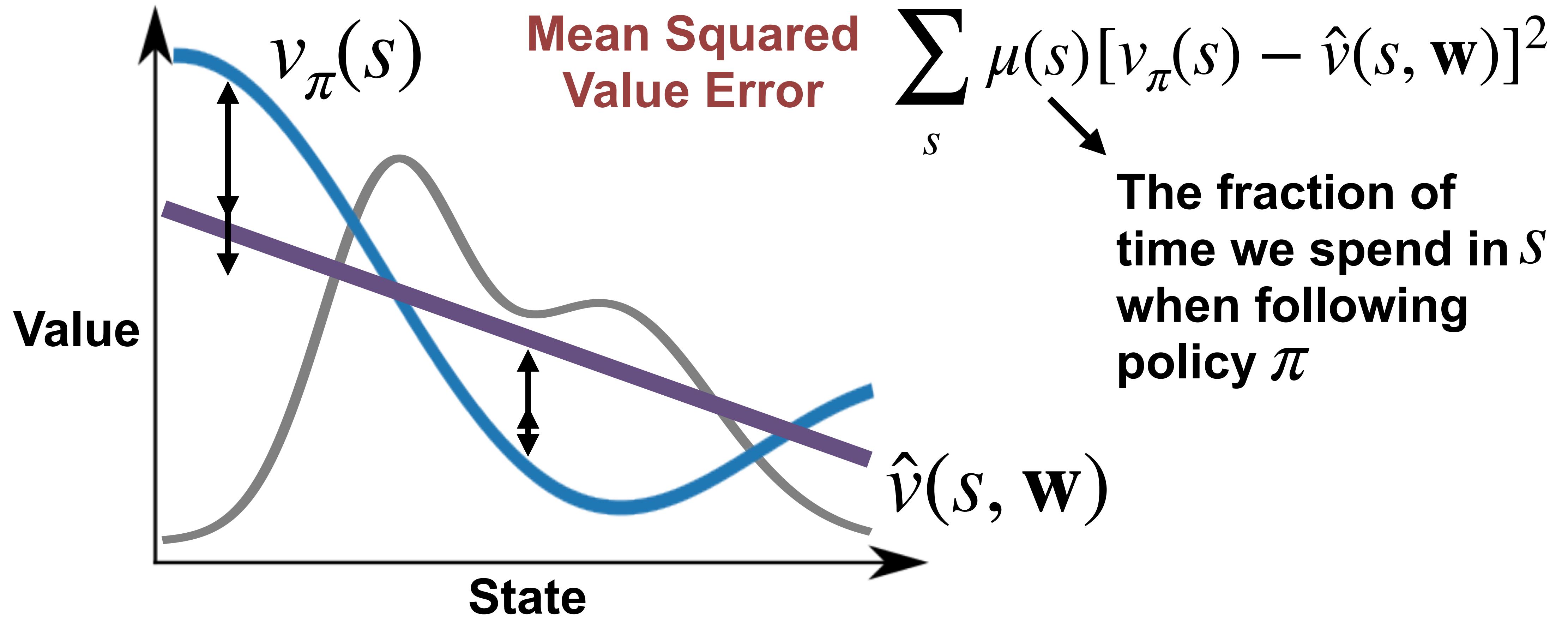
Video 3: Framing Value Estimation as Supervised Learning

- If we can setup the problem of learning a value function (policy evaluation) as a supervised learning problem, then we can borrow methods from supervised learning to do reinforcement learning with function approximation.
- Goals:
 - Understand how value estimation can be framed as a supervised learning problem
 - Recognize that not all function approximation methods are well suited for reinforcement learning.

Video 4: Value Error

- We want to change the parameters of our function to estimate the value. We need an objective function!
- Goals:
 - Understand the mean-squared value error objective for policy evaluation
 - Explain the role of the state distribution in the objective

The Mean Squared Value Error Objective



Question: Why didn't we use the Value Error in the tabular setting?

Video 5: Introducing Gradient Descent

- An algorithm for adapting the parameters of our estimate of the value function.
- Goals:
 - Understand the idea of gradient descent
 - Understand that gradient descent converges to stationary points

Video 6: Gradient Monte Carlo for Policy Evaluation

- We use gradient descent idea to get an online algorithm to adjust the parameters of our value function estimate
- Goals:
 - Understand how to use gradient descent and stochastic gradient descent to minimize value error
 - Outline the gradient Monte Carlo algorithm for value estimation

Video 7: State Aggregation with Monte Carlo

- So far we have said the value function could be any parametric function. Here we use a particular one---state aggregation. Simple and effective. And we run an experiment on a big Random Walk Problem
- Goals:
 - Understand how state aggregation can be used to approximate the value function
 - Apply Gradient Monte-Carlo with state aggregation

Function approximation is something we want, not something we have to do!

- Consider the 1000 state chain:
 - We could store this in a table, but aggregating states makes learning faster and the final error is pretty good!
- Consider GO:
 - The state-space is finite!
 - But the agent will never, ever see all the states! We need to generalize to do well

Video 8: Semi-gradient TD for Policy Evaluation

- TD with function approximation. Now we can learn value functions, in continuous state spaces AND update the value function parameters on every time-step!!
- Goals:
 - Understand the TD-update for function approximation
 - Outline the Semi-gradient TD algorithm for value estimation.

Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$

Input: the policy π to be evaluated

Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$

Algorithm parameter: step size $\alpha > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose $A \sim \pi(\cdot | S)$

 Take action A , observe R, S'

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})$$

$S \leftarrow S'$

 until S is terminal

Question: What is different compared to Tabular TD(0)?

Video 9: Comparing TD and MC with State Aggregation

- An experiment comparing TD and MC with a simple function approximation.
- Goals:
 - Understand that TD converges to biased value estimates
 - Understand that TD can learn faster than Gradient Monte Carlo.

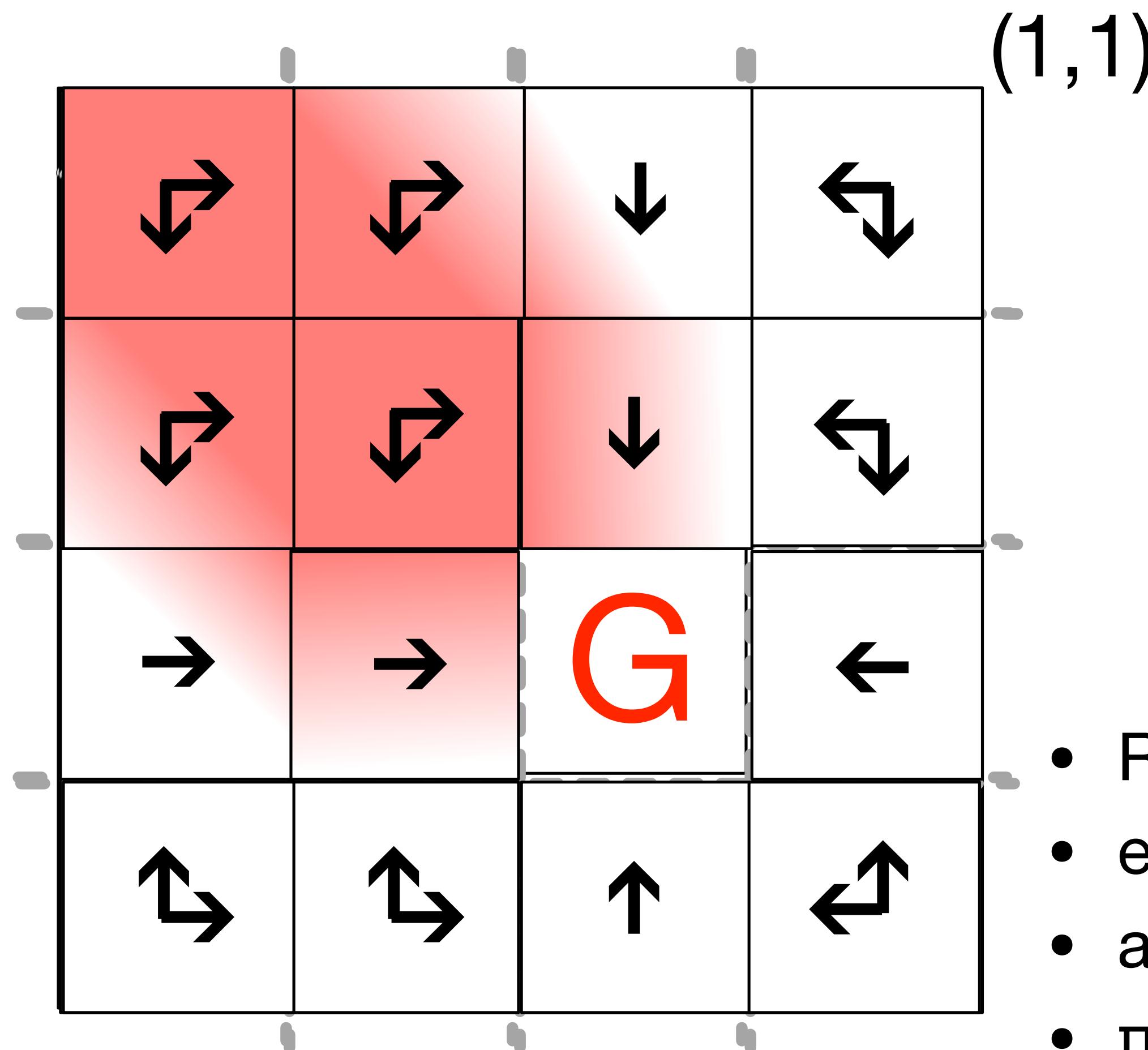
Video 10: The Linear TD Algorithm

- Linear function functions are special. Most of the theory in RL is for the case of linear function approximation. The algorithms can work well, if we have good features.
- Goals:
 - Derive the TD-update with linear function approximation
 - Understand that tabular TD is a special case of linear semi-gradient TD
 - Understand why we care about linear TD as a special case.

Video 11: The True Objective for TD

- A bit of theory about TD with function approximation. What does the algorithm converge to?
- Goals:
 - Understand the fixed point of linear TD
 - Describe a theoretical guarantee on the mean squared value error at the TD fixed point

What might the μ (proportion of time the agent spends in each state) look like with this state aggregation?



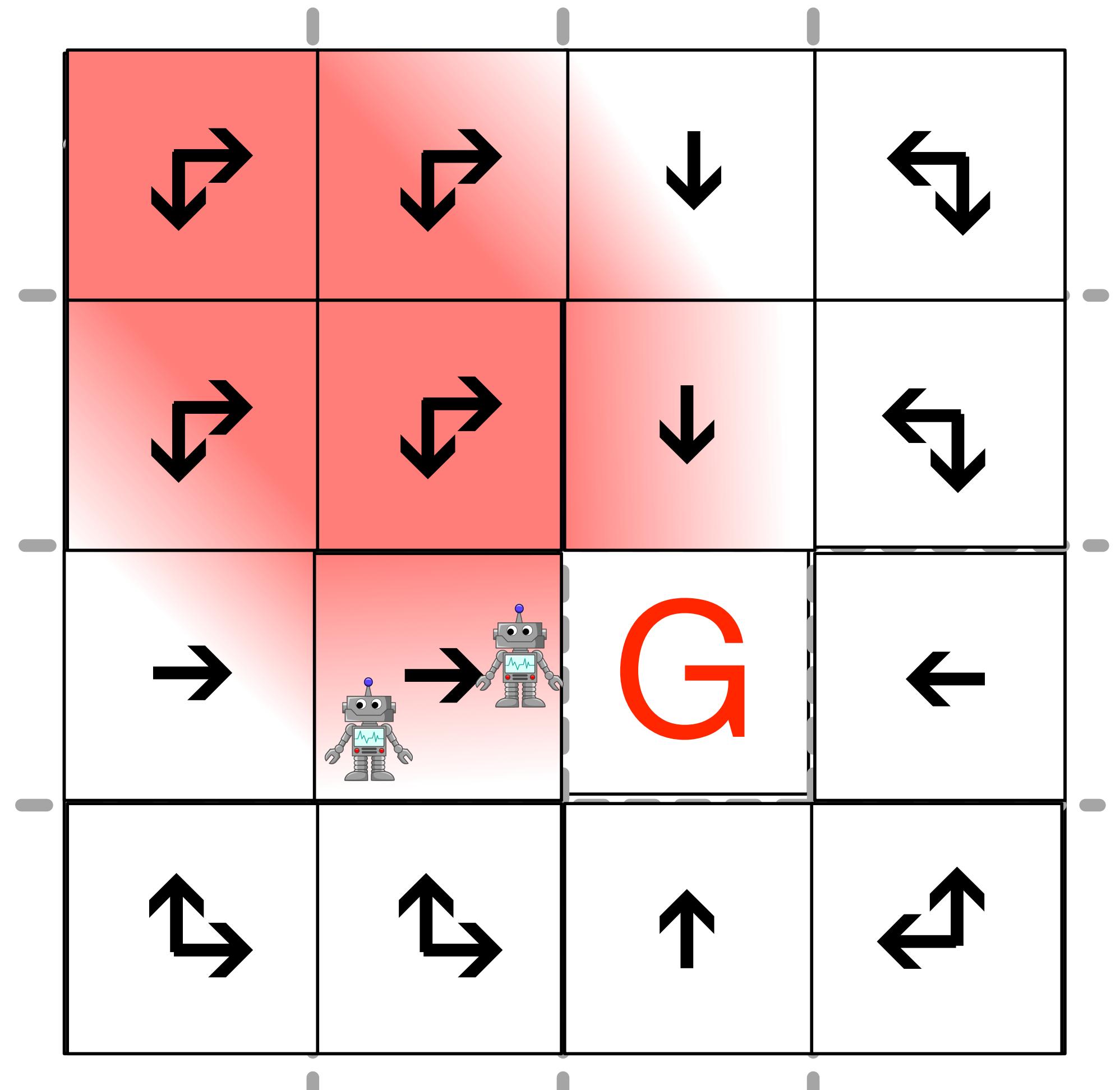
- $R = +1$ per step
- episodic, $\gamma = 1$
- agent starts in the top-left corner
- π = shortest path policy

Mean Squared Value Error

$$\sum_s \mu(s)[v_\pi(s) - \hat{v}(s, w)]^2$$

The fraction of time we spend in s when following policy π

$\mu(s)$ Impacts how we update $\hat{v}(s, w)$



**Mean Squared
Value Error**

$$\sum_s \mu(s)[v_\pi(s) - \hat{v}(s, w)]^2$$

The fraction of time we spend in s when following policy π

The usual recipe for gradient descent

1. Specify a function approximation architecture (parametric form of value function)
2. Write down your objective function
3. Take the derivative of objective function with respect to the weights
4. Simplify general gradient expression for your parametric form
5. Make a weight update rule:
 - $\mathbf{w} = \mathbf{w} - \alpha \text{ gradient}$

gradient descent

lets try out the recipe

1. Specify a function approximation architecture (parametric form of value function)

- We will use **State Aggregation**
 - so the **features** are always **binary** with only a single active feature that is not zero
 - the value function is a **linear function**
 - that is, we query the value function by a simple procedure:
 1. **query the features** for the current state
 2. take the inner product between the features and the weights

$$v_\pi(s) \approx \hat{v}(s, \mathbf{w}) \doteq \mathbf{w}^\top \mathbf{x}(s) \doteq \sum_{i=1}^n w_i \cdot x_i(s)$$

2. Write down your objective function

- We will use the value error

$$\begin{aligned}\overline{VE}(\mathbf{w}) &\doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \boxed{\mathbf{w}^T \mathbf{x}(s)}]^2\end{aligned}$$

state
aggregation

3. Take the gradient of objective function with respect to the weights

$$\nabla \overline{VE}(\mathbf{w}) = \nabla \sum_{s \in \mathcal{S}} \mu(s) [\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)]^2$$

3. Take the gradient of objective function with respect to the weights

$$\begin{aligned}\nabla \overline{VE}(\mathbf{w}) &= \nabla \sum_{s \in \mathcal{S}} \mu(s) [\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)]^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) \nabla [\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)]^2 \\ &= - \sum_{s \in \mathcal{S}} \mu(s) 2[\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \nabla \mathbf{w}^T \mathbf{x}(s)\end{aligned}$$

4. Simplify the general gradient expression to be specific for your parametric form

$$\nabla_{\mathbf{w}}^T \mathbf{x}(s) = \mathbf{x}(s)$$

The gradient of the inner product is just \mathbf{x}

gradient of

linear value function

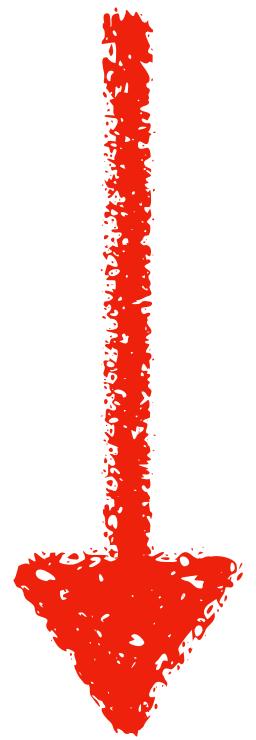
approximation (state agg.)

4. Simplify general gradient ...

$$\begin{aligned}\nabla \overline{VE}(\mathbf{w}) &= \nabla \sum_{s \in \mathcal{S}} \mu(s) [\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)]^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) \nabla [\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)]^2 \\ &= - \sum_{s \in \mathcal{S}} \mu(s) 2[\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \nabla \mathbf{w}^T \mathbf{x}(s) \\ &= - \sum_{s \in \mathcal{S}} \mu(s) 2[\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \boxed{\mathbf{x}(s)}\end{aligned}$$

5. Make weight update rule: $\mathbf{w} = \mathbf{w} \boxed{-} \alpha$ gradient

$$\nabla \overline{VE}(\mathbf{w}) = \boxed{-} \sum_{s \in \mathcal{S}} \mu(s) 2[\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \mathbf{x}(s)$$



$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha 2[\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \mathbf{x}(s)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [\nu_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \mathbf{x}(s)$$

Wait, Wait, Wait!! We don't have v_π

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [v_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \mathbf{x}(s)$$

Let's replace it with something we do have!

Let's replace v_π with something we do have!

Let's call it's replacement U_t

Whatever we use in place of v_π , it should satisfy one criteria!

$$v_\pi(s) = \mathbb{E}[U_t]$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha [v_\pi(s) - \mathbf{w}^T \mathbf{x}(s)] \mathbf{x}(s)$$

Whatever we use in place of v_π , it should satisfy one criteria!

$$v_\pi(s) = \mathbb{E}[U_t]$$

We know one such replacement, that meets this criteria!

$$U_t \doteq G_t$$

A sample of the return!!

Since we are using sample returns we have a Monte Carlo algorithm!

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha[G_t - \mathbf{w}^T \mathbf{x}(s)]\mathbf{x}(s)$$

Monte Carlo Policy Evaluation for finding v_π

Exercise Question

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_s \mu(s) [\nu_\pi(s) - \hat{\nu}(s, \mathbf{w})]^2$$

- Why can't we directly optimize the MSVE? We know the stochastic gradient descent update would be the following

$$\mathbf{w}_t + \alpha [\nu_\pi(S_t) - \hat{\nu}(S_t, \mathbf{w})] \nabla \hat{\nu}(S_t, \mathbf{w})$$

- Further, why doesn't the TD fixed point minimize the MSVE?