

# Advanced Topics on Artificial Intelligence

Alban Grastien

The Australian National University

Second Semester, 2020

# Model-free decision

- We are now considering the problem of decision under uncertainty when **the model is unknown**
  - We do not know the transition function  $P(s, a, s')$
  - We do not know the reward function  $R(s, a, s')$
  - But the state is still fully observable

# My Previous Lecture: Bandit

Two important problems:

- How to update the knowledge
  - E.g., average outcome, temporal difference
- How to balance exploration / exploitation
  - E.g.,  $\epsilon$ -greedy strategy or UCB

# Today

- Evaluation of a policy
- Computing the optimal (?) policy: SARSA / Q-learning

# Yesterday / today

Context: Reinforcement Learning (the  $P$  and  $R$  functions of the MDP are unknown)

- Yesterday:
  - How to estimate the value of a policy
  - How to compute a good / optimal policy anytime (SARSA, Q-Learning)
- Today: Problems with large number of states
  - Parametrised policies
  - Finding good parametrised policies:
    - Gradient descent
    - Analytical solutions

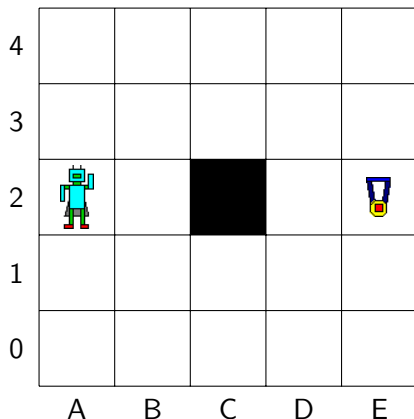
# Issue with SARSA/Q-Learning

- If there are too many states (e.g.,  $10^9$  and more), it is impossible to maintain an estimate of  $V(s)$  or  $Q(s, a)$  for each state/action
- These methods are also not able to generalise:
  - if two states are very similar, then the optimal action is very likely to be the same in both states

# Parametrised Policies

- A **parametrised policy** is a function  $\pi_{\theta}(s) \rightarrow \text{Prob}(A(s))$  where  $\theta$  is a vector (list of parameters)
  - Different values to the parameters yield a different policy
- We write:  $\theta = [\theta_1, \theta_2, \dots, \theta_n]$
- Next slides: examples

# Example of Parametrised Policy (1)



- We compute a score  $f(s, a)$  for each state  $s$  and each action  $a$ :  

$$f((x, y), \uparrow) = \min(\varepsilon, x\theta_1 + y\theta_2)$$

$$f((x, y), \downarrow) = \min(\varepsilon, x\theta_3 + y\theta_4)$$

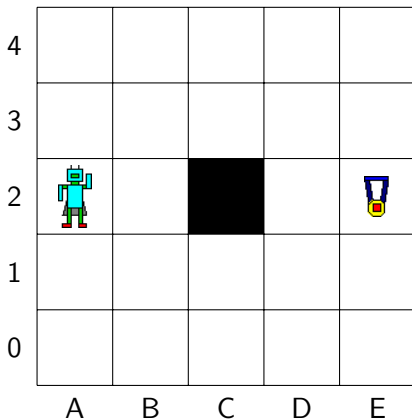
$$f((x, y), \rightarrow) = \min(\varepsilon, x\theta_5 + y\theta_6 y)$$

$$f((x, y), \leftarrow) = \min(\varepsilon, x\theta_7 + y\theta_8 y)$$
- Here  $\theta = [\theta_1, \dots, \theta_8]$
- If  $a \in A(s)$ , the probability of applying action  $a$  from  $s$  is:  

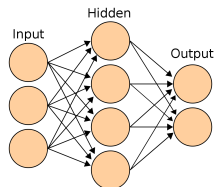
$$\pi_{\theta}(s)(a) = \frac{f(s, a)}{\sum_{a' \in A(s)} f(s, a')}$$



## Example of Parametrised Policy (2)



### Neural network



- Parameters = weights of the connections between the neurons
- What are the inputs of the NN?

Image: By en>User:Cburnett - Own work

This W3C-unspecified vector image was created with Inkscape.,  
CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=1496812>

# Optimal Parameters

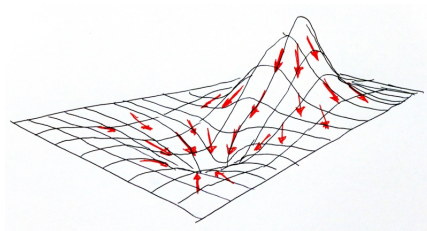
- Given the parameters  $\theta$ , the value of the initial state  $s_0$  for  $\pi_\theta$  is  $V_\theta(s_0) := V_{\pi_\theta}(s_0)$

- The optimal parameters  $\theta^*$  are those such that

$$V_{\theta^*}(s_0) \geq V_\theta(s_0) \quad \text{for all parameters } \theta$$

- There may be no parameter that is optimal **for all states**
- From now on, I will use the notation  $V(\theta)$  to represent  $V_\theta(s_0)$

# Gradient Descent — General Idea



- The derivative  $\frac{\partial f}{\partial x}(x)$  of function  $f(x)$  indicates in which direction to go in order to increase / decrease  $f(x)$
- The gradient is similar for several parameters

# Gradient Descent

- Assume that  $\pi_{\theta}(s)$  is differentiable wrt  $\theta$
- The **gradient** of the  $V(\theta)$  is the vector:

$$\nabla V(\theta) = \left[ \frac{\partial V(\theta)}{\partial \theta_1}, \frac{\partial V(\theta)}{\partial \theta_2}, \dots, \frac{\partial V(\theta)}{\partial \theta_n} \right]$$

- Let  $\alpha$  be a **step size**.
- Then we update the parameters  $\theta' := \theta + \alpha \nabla V(\theta)$
- If  $\alpha$  is small enough, then  $V(\theta')$  is better than  $V(\theta)$

# How to Compute the Gradient?

First method: Estimate

- Let  $\mathbf{1}_i$  be the vector  $\underbrace{[0, \dots, 0, 1, 0, \dots, 0]}_{i-1}$

- By definition:

$$\frac{\partial V(\boldsymbol{\theta})}{\partial \theta_i} = \lim_{\varepsilon \rightarrow 0} \frac{V(\boldsymbol{\theta}') - V(\boldsymbol{\theta})}{\varepsilon}$$

where  $\boldsymbol{\theta}' := \boldsymbol{\theta} + \varepsilon \mathbf{1}_i$  is a perturbation of the parameter vector  $\boldsymbol{\theta}$

- Just choose  $\varepsilon$  small enough and assume

$$\frac{\partial V(\boldsymbol{\theta})}{\partial \theta_i} \simeq \frac{V(\boldsymbol{\theta}') - V(\boldsymbol{\theta})}{\varepsilon}$$

# Issues of this Method for Gradient Descent

- Requires a lot of computation, because  $V(\theta)$  needs to be accurate
- Gradient descent often ends up in local optima

# How to Compute the Gradient?

## Second method: Analytically

- Assumptions & Notations
  - Episodic problems (with final states)
  - A run, or a history, is  $h = [s_0, a_1, s_1, a_2, \dots, s_{k_h}]$
  - The set of all histories is  $\mathcal{H}$

- The value associated with parameters  $\theta$  is

$$V(\theta) = \sum_{h \in \mathcal{H}} P(h \mid \theta) \times R(h).$$

- We want to compute

$$\nabla_{\theta} V(\theta)$$

- Notice that  $P(h \mid \theta)$  is unknown (and  $R(h)$  too)!

# Analytical method (cont.)

$$\nabla_{\boldsymbol{\theta}} V(\boldsymbol{\theta})$$

- $= \nabla_{\boldsymbol{\theta}} \sum_{h \in \mathcal{H}} \left( P(h | \boldsymbol{\theta}) \times R(h) \right)$
- $= \sum_{h \in \mathcal{H}} \left( R(h) \times \nabla_{\boldsymbol{\theta}} P(h | \boldsymbol{\theta}) \right)$
- $= \sum_{h \in \mathcal{H}} \left( R(h) \times \frac{P(h|\boldsymbol{\theta}) \times \nabla_{\boldsymbol{\theta}} P(h|\boldsymbol{\theta})}{P(h|\boldsymbol{\theta})} \right)$
- $= \sum_{h \in \mathcal{H}} \left( R(h) \times P(h | \boldsymbol{\theta}) \times \nabla_{\boldsymbol{\theta}} \log(P(h | \boldsymbol{\theta})) \right)$
- $\simeq \frac{1}{m} \sum_{i \in \{1, \dots, m\}} \left( R(h^i) \times \nabla_{\boldsymbol{\theta}} \log(P(h^i | \boldsymbol{\theta})) \right)$  where  $h^1, \dots, h^m$  are sampled trajectories



# Analytical method (cont.)

Looking at a single history  $h = [s_0, a_1, s_1, a_2, \dots, s_k]$

$$\nabla_{\theta} \log(P(h \mid \theta))$$

$$\bullet = \nabla_{\theta} \log(P(s_0 \mid \theta) \times P(a_1 \mid s_0, \theta) \times P(s_1 \mid s_0, a_1, \theta) \times P(a_2 \mid s_1, \theta) \dots)$$

$$\bullet = \nabla_{\theta} \left( \log(P(s_0 \mid \theta)) + \log(P(a_1 \mid s_0, \theta)) + \log(P(s_1 \mid s_0, a_1, \theta)) + \log(P(a_2 \mid s_1, \theta)) + \dots \right)$$

$$\bullet = \nabla_{\theta} \log(P(s_0 \mid \theta)) + \nabla_{\theta} \log(P(a_1 \mid s_0, \theta)) + \nabla_{\theta} \log(P(s_1 \mid s_0, a_1, \theta)) + \nabla_{\theta} \log(P(a_2 \mid s_1, \theta)) \dots$$

$$\bullet = 0 + \nabla_{\theta} \log(P(a_1 \mid s_0, \theta)) + 0 + \nabla_{\theta} \log(P(a_2 \mid s_1, \theta)) + \dots$$

This value can be computed because we know  $P(a \mid s, \theta)$  for each  $s, a, \theta$

# AlphaGo

Silver, David, et al. "**Mastering the game of go without human knowledge.**" *Nature* 550.7676 (2017): 354-359.

Three ingredients:

- Monte Carlo Tree Search
- A neural network that computes good moves (used for simulation by MCTS)
- A neural network that predicts the winner of a game (to prune the search depth)

# Last Word

[https://www.youtube.com/watch?v=\\_cQITY0SPiw](https://www.youtube.com/watch?v=_cQITY0SPiw)

Discussion about past, present, and future of AI