

Reinforcement Learning in Finance

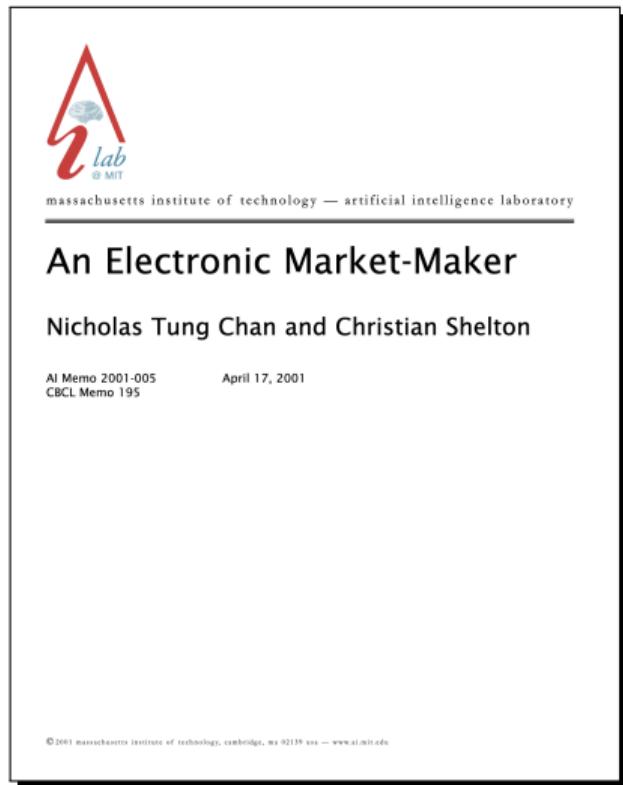
Leveraging Structure (for Fun and Profit)

Thomas Spooner

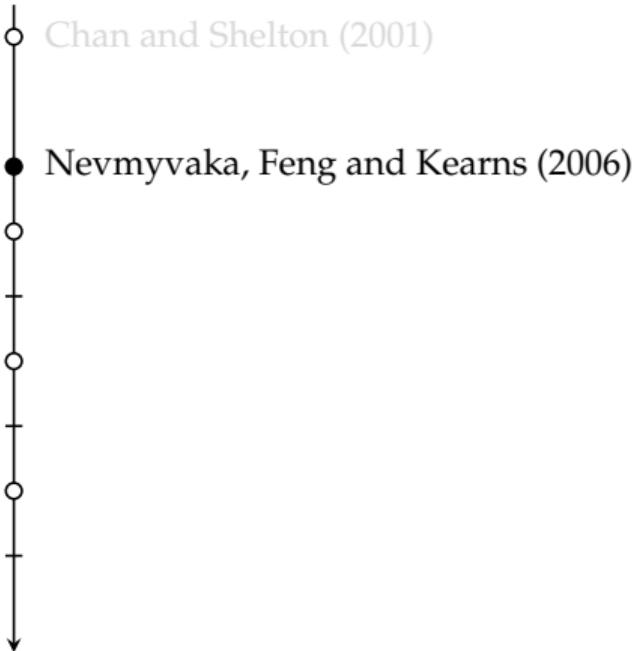
13th August, 2022

Key Papers: A Timeline

- Chan and Shelton (2001)
-
-
- +
-
- +
-
- +
- ↓



Key Papers: A Timeline



Reinforcement Learning for Optimized Trade Execution

Yuriy Nevmyvaka

Lehman Brothers, 745 Seventh Av., New York, NY 10019, USA

Yi Feng

Michael Kearns[†]

University of Pennsylvania, Philadelphia, PA 19104, USA

yuriy.nevmyvaka@lehman.com

fengyi@cis.upenn.edu

mkearns@cis.upenn.edu

Abstract

We present the first large-scale empirical application of reinforcement learning to the important problem of optimized trade execution in modern financial markets. Our experiments are based on 1.5 years of millisecond time-scale limit order data from NASDAQ and demonstrate the promise of reinforcement learning methods to market microstructure problems. Our learning algorithm introduces and exploits a natural “low-impact” factorization of the state space.

1. Introduction

In the domain of quantitative finance, there are many decision problems in which there are precisely specified objectives and constraints, but no clear way to solve the problem of optimized trade execution (which has also been examined in the theoretical computer science literature as the one-way trading problem). Previous work on this problem includes (Argen and Chiaro, 2002), (Brennan et al., 2003), (Huang et al., 2003), (Huang et al., 2007), (Huang et al., 2009), (Yann et al., 2001), (Kakade et al., 2004), and (Nevmyvaka et al., 2005). In this problem, the goal is to sell (respectively, buy) V shares of a given stock within a fixed time period (or horizon) T , while trying to maximize the revenue (and, respectively, minimize the capital spent). This problem arises frequently in practice as the result of “higher level” investment decisions. For example, if a trading strategy, whether automated or not, by some means decides that Qualcomm stock is likely to decrease in value, it may decide to short-sell a block of shares. But this ‘decision’ is still underspecified, since various trade-offs may arise, such as between the speed of execution and the cost of execution (see Section 2). The problem of optimized execution can thus be viewed as an important horizontal aspect of quantitative trading, since virtually

every strategy’s profitability will depend on how well it is implemented.

For most of the history of the U.S. equity markets, the only information available to an agent attempting to optimize trade execution was the sequence of prices of stock-exchange trades over the current bid or ask (the last outstanding buy and sell orders, respectively). But recent years, electronic markets such as NASDAQ have begun releasing, in real time, all of the outstanding buy and sell limit order prices and volumes (often referred to as order books). This is a remarkable development, because it is clear that such data might be extremely valuable for the problem of optimized execution, since the distribution of outstanding limit orders may help predict short-term price direction, likelihood of execution at a given price, buy or sell side imbalances, and so on.

In this paper, we present the first extensive empirical application of reinforcement learning (RL) to the problem of optimized execution using large-scale NASDAQ market microstructure data sets. Our basic experimental methodology (detailed in Section 3) consists of the following steps:

1. The construction of a set of (observed) state variables that can be derived from the microstructure data and the RL algorithm’s own actions. These variables are divided into “private” variables (such as the amount of time and shares remaining for the algorithm in the execution problem) and “market” variables (which reflect various features of the trading activity in the stock).

2. The application of a collection of RL algorithms, which consist of strong baselines for the execution problem to improve computational efficiency, to extremely large datasets. Our empirical results are based on 1.5 years of very high-frequency (millisecond time scale) microstructure data for several NASDAQ stocks.

3. The empirical comparison of the RL-optimized execution policy to a variety of natural baseline execution strategies.

Portions of this work were conducted while the authors were in the Equity Strategies department of Lehman Brothers in New York City. Appearing in Proceedings of the 23rd International Conference on Machine Learning, Pittsburgh, PA, 2006. Copyright 2006 by the author(s)/owner(s).

Key Papers: A Timeline

- Chan and Shelton (2001)
- Nevmyvaka, Feng and Kearns (2006)
- Avellaneda and Stoikov (2008)

Quantitative Finance, Vol. 8, No. 3, April 2008, 217–224
Downloaded By [NIH Statal School of Medicine, Levy Library] At: 23:19 28 April 2008

Routledge
Taylor & Francis Group



High-frequency trading in a limit order book

MARCO AVELLANEDA and SASHA STOIKOV*

Mathematics, New York University, 251 Mercer Street, New York, NY 10012, USA

(Received 24 April 2006; in final form 3 April 2007)

1. Introduction

The role of a dealer in securities markets is to provide liquidity on the exchange by quoting bid and ask prices at which he is willing to buy and sell a specific quantity of assets. Traditionally, this role has been filled by market-maker or specialist firms. In recent years, with the growth of electronic exchanges such as Nasdaq's Insti, anyone willing to submit limit orders in the system can effectively play the role of a dealer. Therefore, the rise of high-frequency data on the limit order book (see www.Quote.com) ensures a fair playing field where various agents can post limit orders at the prices they choose. In this paper, we study the optimal submission strategies of bid and ask orders in such a limit order book.

The pricing strategies of dealers have been studied extensively in the microstructure literature. The two most often addressed sources of risk facing the dealer are (i) the

inventory risk arising from uncertainty in the asset's value and (ii) the asymmetric information risk arising from informed traders. Useful surveys of their results can be found in Blaiss *et al.* (2004), Stoll (2003) and a book by O'Hara (1997). In this paper, we will focus on the inventory effect. In fact, our model is closely related to a paper by Ho and Stoll (1981), which analyses the optimal prices for a monopolistic dealer in a single stock. In their model, the dealer submits limit orders for the asset, and derive optimal bid and ask quotes around this price, to account for the effect of the inventory. This inventory effect was found to be significant in an empirical study of AMEX Options by Ho and Macris (1984). In another paper by Ho and Stoll (1989), the problem of dealers under competition is analysed and the bid and ask prices are shown to be related to the reserves (or indifference) prices of the dealers. In our framework, we assume that our agent is the one player in the market and the 'true' price is given by the market mid-price.

Of crucial importance to us will be the arrival rate of buy and sell orders that will reach our agent. In order

*Corresponding author. Email: usashstoikov@gmail.com

Key Papers: A Timeline



Adaptive Market Making via Online Learning

Jared Abernethy*
Computer Science and Engineering
University of Michigan
jabernet@umich.edu

Satyen Kale
IBM T. J. Watson Research Center
skale@us.ibm.com

Abstract

We consider the design of strategies for market making in an exchange. A market maker generally seeks to profit from the difference between the buy and sell price of an asset, yet the market maker also takes exposure risk in the event of large price movements. Profit guarantees for market making strategies have typically required certain stochastic assumptions on the price fluctuations of the asset in question; for example, assuming a model in which the price process is mean reverting. We prove that such guarantees can be obtained even under worst case (adversarial) settings. We prove structural properties of these strategies which allows us to design a master algorithm which obtains low regret relative to the best such strategy in hindsight. We run a set of experiments showing favorable performance on recent real-world stock price data.

1 Introduction

When a trader enters a market, say a stock or commodity market, with the desire to buy or sell a certain quantity of an asset, how is this trader guaranteed to find a counterparty to engage in a transaction at a reasonable price? This is not a problem in a liquid market, with a deep pool of traders ready to buy or sell at any time, but in a thin market the lack of counterparties can be troublesome. A naked trader may even be willing to transact at a worse price in exchange for immediate execution.

This is where a *market maker* (MM) can be quite useful. A MM is any agent that participates in a market by offering to both buy and sell the underlying asset at any time. To put it simply, a MM covers the guarantee of liquidity to the marketplace by providing the ability to trade for any trader. The act of market making has two primary benefits and risks. On one side, the MM bears the risk of transacting with better-informed traders that may know much more about the movement of the asset's price, and in such scenarios the MM can take on a large inventory of shares that it may have to offload at a worse price. On the positive side, the MM can profit as a result of the bid-ask spread, the difference between the buy price and the sell price of the asset. For example, if the MM buys 100 shares of a stock from one trader at a price of p , and then immediately sells 100 shares of stock to another trader at a price of $p + \Delta$, the MM records a profit of 100Δ .

This describes the central goal of a profitable market making strategy: minimize the inventory risk of large movements in the price while simultaneously aiming to benefit from the bid-ask spread. The MM strategy has a state, which is the current inventory or holdings, receives as input order and price information from the market, and queries what price to offer in the market. In the present paper we assume that the MM interacts with a constant double auction via an order book, and the MM can place both market and limit orders to the order book.

A number of MM strategies have been proposed, and in many cases certain profitloss guarantees have been given. But to the best of our knowledge all such guarantees (aside from [4]) have required

*Work performed while the author was in the CIS Department at the University of Pennsylvania and funded by a Simons Postdoctoral Fellowship

Key Papers: A Timeline

- 
- Chan and Shelton (2001)
 - Nevmyvaka, Feng and Kearns (2006)
 - Avellaneda and Stoikov (2008)
 - Abernethy and Kale (2013)
 - Spooner et al. (2018)

Market Making via Reinforcement Learning

Thomas Spooner
Department of Computer Science
University of Liverpool
t.spooner@liverpool.ac.uk

Rahul Savani
Department of Computer Science
University of Liverpool
rahul.savani@liverpool.ac.uk

arXiv:1804.04216v1 [cs.AI] 11 Apr 2018

ABSTRACT

Market making is a fundamental trading problem in which an agent provides liquidity by continually offering to buy and sell a stock. This is a challenging decision-making task due to the need of accelerating an unobservable execution and ultimately losing money. In this paper, we develop a high-fidelity simulation of limit order book markets, and use it to design a market making agent using temporally-difference reinforcement learning. We use a linear model of tick evolution as the environment, and propose and design a custom reward function that controls inventory risk. We demonstrate the effectiveness of our approach by showing that our agent outperforms both simple benchmark strategies and a recent online learning approach from the literature.

KEYWORDS
Market Making, Limit Order Books, TD-Learning, TDE Coding

1 INTRODUCTION

The role of a market maker is to provide liquidity by facilitating transactions with other market participants. Like reading trading news or analysing financial reports, market making requires the use of the electronic limit order book (LOB), as the need to handle more data and act over shorter time scales renders the task almost impossible for humans [25, 29]. Upwards of 60% of trading volume on most exchanges is now done by automated systems, and thus it is important to understand how these systems can be used to automatically trade stocks [27, 36]. This paper uses reinforcement learning (RL) to design competitive market making agents for financial markets using high-frequency historical equities data.

1.1 Related work

Market making has been studied across a number of disciplines, including economics, finance, artificial intelligence (AI), and machine learning. A classic approach in the finance literature is to treat market making as a problem of stochastic optimal control. Here, a model for the market dynamics is built, and the agent's actions and the rewards for the resulting dynamics are designed [5, 10, 11, 20, 23, 25]. Recent results in this line of research have studied price impact, adverse selection and predictability [1], and augmented the problem characteristics with risk measures and inventory constraints [7, 22].

Proc. of the 17th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '18). International Foundation for Autonomous Agents and Multiagent Systems (www.ijamas.org). All rights reserved.

John Fearnley
Department of Computer Science
University of Liverpool
fearnley@liverpool.ac.uk

Andreas Koukoulinis
Stratagene Technologies Ltd,
and University College London
andreas@stratagene.co

Another prominent approach to studying market making and limit order book markets has been that of zero-intelligence (ZI) agents. A study of 21 agents has been conducted by Abernethy and Kale [1]. These agents do not require training, research, or memory, but can, for example, adhere to inventory constraints [18]. Never, more intelligent variants, now even incorporate learning mechanisms [14, 42]. Here, agents are explicitly evaluated in simulated markets without trading against ZI agents.

A significant body of literature, in particular in AI, has studied the market making problem for prediction markets [6, 34, 35]. In this setting, the agent's main goal is to elicit information from informed participants in the market. While later studies have addressed profitability, the setup remains quite distinct from the financial one considered here.

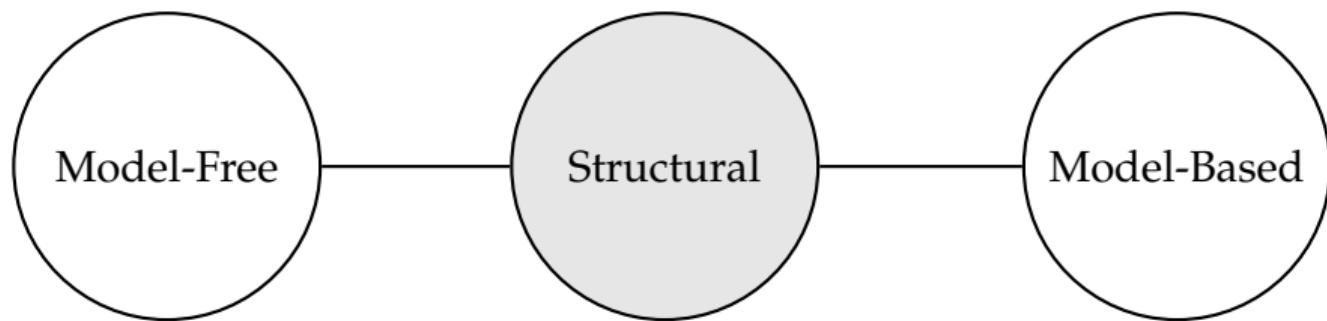
Reinforcement learning has been applied for other financial trading problems [12, 37, 39], including order execution [31] and financial market prediction [38]. The first case of applying RL to market making [12] focused on the impact of market size (i.e., to instrumented traders) on the agent's quoting behaviour and showed that RL successfully converges on the expected strategies for a number of centralized environments. They did not however capture the effect of market dynamics with complex bidding, capturing bid and cancellation, nor the complexities of using continuous state variables. Moreover, [12] found that temporal-difference RL struggled in their setting, a finding echoed in [39]. [12] attributed this to parameter sensitivity and the lack of a clear baseline for comparison, despite the relative simplicity of their market simulation. In follow up work, [38] and [39] propose sampling as a solution to the problems observed with off-policy learning. In contrast, we find temporal-difference RL to be effective for the market making problem, and we also show that it is able to directly steer our function approximator and reward function.

One of the most recent related works is [21], which uses an online learning approach to develop a market making agent. They prove that their agent is able to learn to make profits and to stabilize their orders under strong assumptions on executions. For example, they assume that the market has sufficient liquidity to execute market orders entirely at the posted price with no slippage. We use this approach as one of the benchmarks for our empirical evaluation and address the impact of trading in a more realistic environment.

What is Structure?



What is Structure?



Reinforcement Learning

Algorithmic Trading

Leveraging Structure: *Gradient Factorisation*

Leveraging Structure: *Transition Acyclicity*

Reinforcement Learning

Algorithmic Trading

Gradient Factorisation

Transition Acyclicity

Reinforcement Learning

Markov Chains

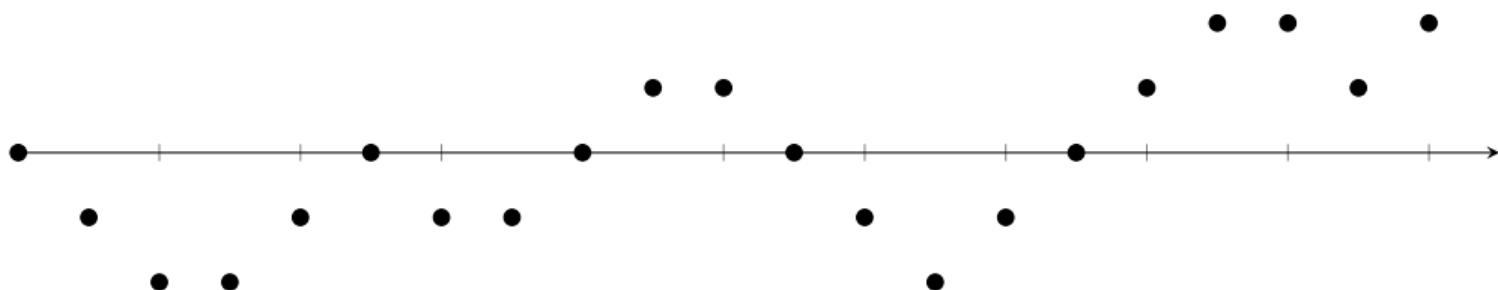
A *Markov chain* (MC) $X = (X_0, X_1, \dots)$ is a discrete-time stochastic process.

The chain X occupies a *state-space* X s.t. $X_n \in \mathsf{X}$ for all $n \in \mathbb{N}$.

Markov Chains

A *Markov chain* (MC) $X = (X_0, X_1, \dots)$ is a discrete-time stochastic process.

The chain X occupies a *state-space* X s.t. $X_n \in \mathsf{X}$ for all $n \in \mathbb{N}$.

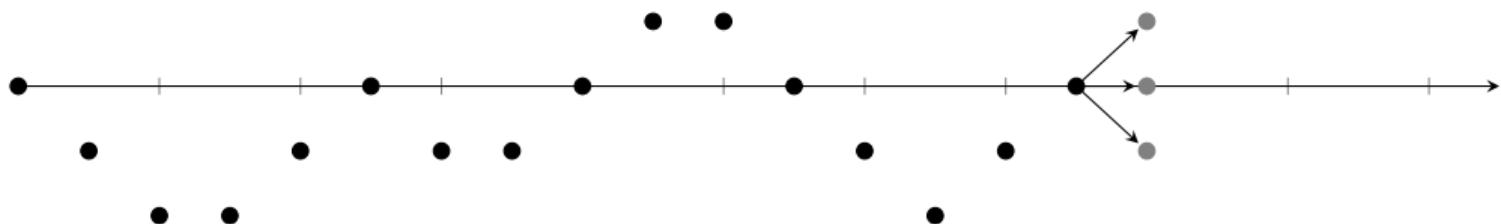


Markov Chains: Dynamics

A chain behaves according to a *transition kernel*:

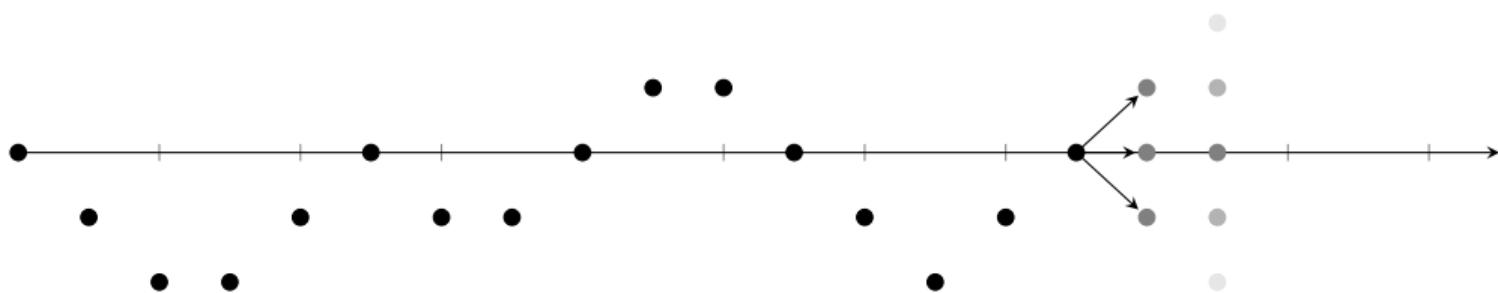
$$\mathbf{P} \doteq \{P(x, B) : x \in \mathsf{X}, B \in \mathfrak{B}(\mathsf{X})\},$$

where $P(x, B) \doteq \mathbb{P}(X_{n+1} \in B \mid X_n = x)$.



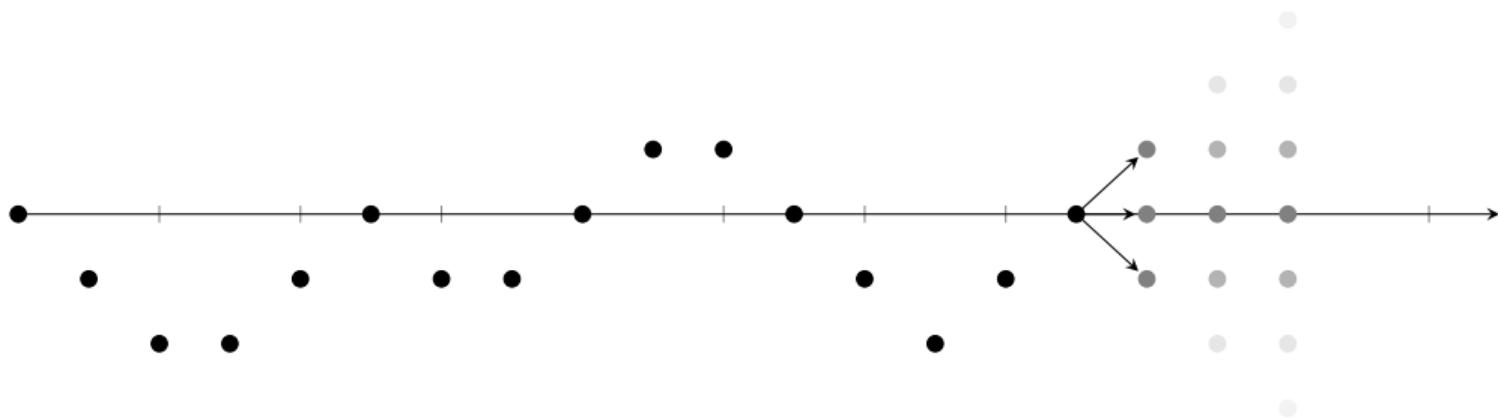
Markov Chains: Dynamics

We can naturally extend to *multi-step kernels*: $P_2(x, B) \doteq \int_X P(x, dy) \cdot P(y, B)$.



Markov Chains: Dynamics

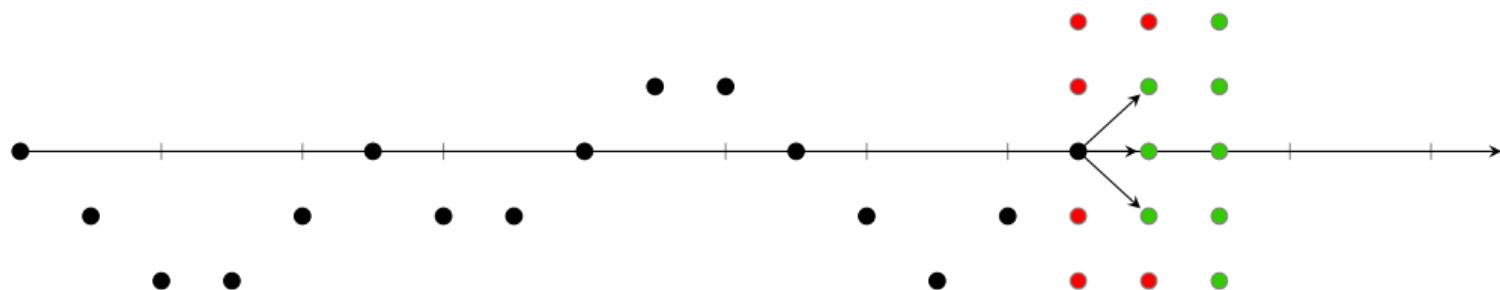
And the rest are defined inductively: $P_n(x, B) \doteq \int_X P(x, dy) \cdot P_{n-1}(y, B)$.



Markov Chains: Reachability

The n -step kernels describe the probability of *reaching* states in X .

A state $x \in \mathsf{X}$ is said to *lead to* a set $B \in \mathfrak{B}(\mathsf{X})$ if $P_n(x, B) > 0$ for some $n > 0$.



Markov Chains: Reachability

What about the probability of *ever returning to a set* $B \in \mathfrak{B}(X)$?

First Return

The time of first return to $B \in \mathfrak{B}(X)$ is

$$\tau_n(B) \doteq \min \{m \geq 1 : X_{n+m} \in B\}.$$

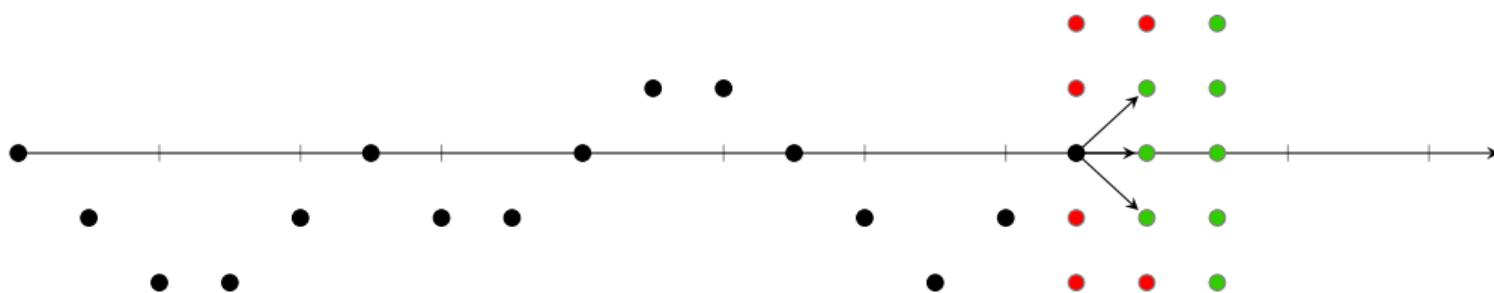
Reachability

The probability of ever reaching B is then

$$P_\infty(x, B) \doteq \mathbb{P}(\tau_n(B) < \infty \mid X_n = x).$$

Markov Chains: Reachability

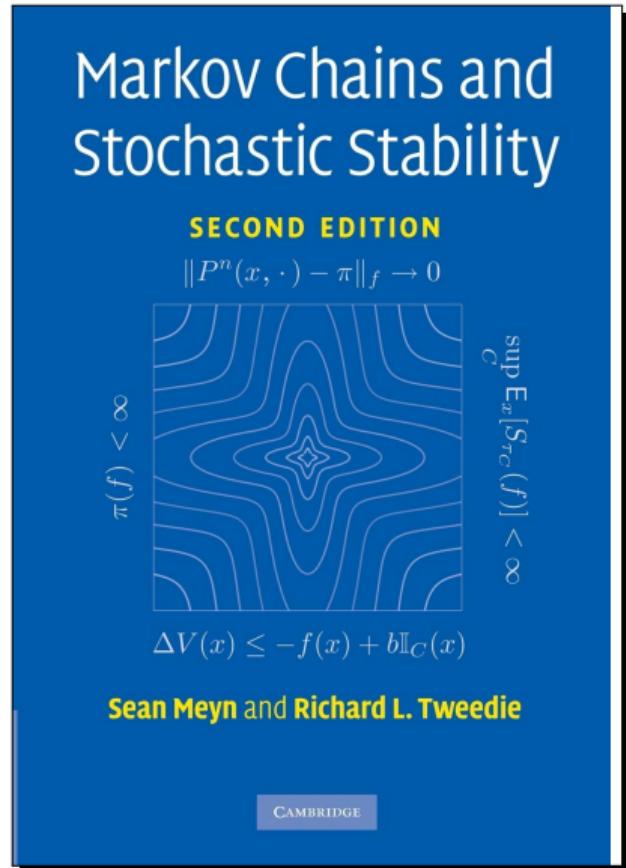
If $P_\infty(x, B) > 0$ then there exists a path from $x \in X$ to $B \in \mathcal{B}(X)$.



Markov Chains

For further reading:

- ▶ Meyn and Tweedie (2012).
- ▶ Grimmett and Stirzaker (2006).



Markov Decision Processes

Markov decision processes (MDPs) extend MCs by the inclusion of:

1. an action space U ; and
2. a reward function $r : X \times U \times X \rightarrow \mathbb{R}$.

The former leads us to define a *joint state-action process* (X, U) , where

$$P_1^\pi(x, B) \doteq \int_{U(x)} P_1(x, B \mid u) \cdot dF_\pi(u \mid x).$$

Markov Decision Processes: Policies

The *behaviour* of an agent is described by a *policy* $\pi \in \Pi$, where

$$\Pi \doteq \{\mathbf{X} \rightarrow \mathcal{P}(\mathbf{U})\}.$$

Deterministic Policies

“If price is higher than c , then buy...”

$$\pi(u | x; c) = \delta_u(\mathbf{1}_{x \geq c})$$

Stochastic Policies

“Quote prices about centre $f_{\theta}(x)\dots$ ”

$$\pi(u | x; \theta) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{(f_{\theta}(x)-x)^2}{2}}$$

Markov Decision Processes: Dynamics

The process (X, U) generates a trajectory according to:

$$X_0 \sim d_0(\cdot),$$

$$U_0 \sim \pi(\cdot | X_0),$$

$$X_1 \sim P_1^\pi(X_0, \cdot),$$

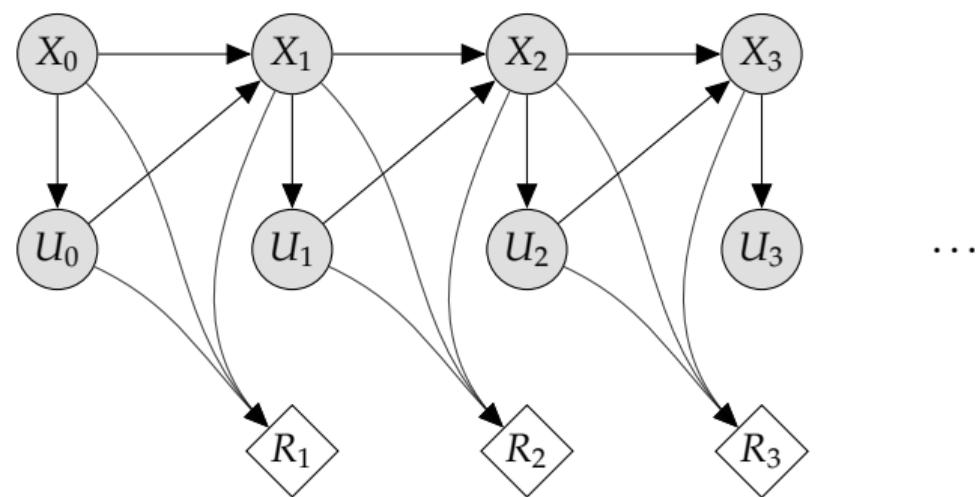
$$R_1 = r(X_0, U_0, X_1),$$

$$U_1 \sim \pi(\cdot | X_1),$$

$$X_2 \sim P_1^\pi(X_1, \cdot),$$

$$R_2 = r(X_1, U_1, X_2),$$

⋮



Markov Decision Processes: Performance Criteria

How do we measure the performance of a policy $\pi \in \Pi$?

$$G_n^\gamma \doteq \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{n+1} + \gamma G_{n+1}.$$

Undiscounted Setting

Let $\gamma \rightarrow 1$ and define

$$J(\pi) = \mathbb{E}_\pi [G_0^1 \mid X_0 \sim d_0]$$

Discounted Setting

Let $\gamma \in [0, 1)$ and define

$$J(\pi) = \mathbb{E}_\pi [G_0^\gamma \mid X_0 \sim d_0]$$

Markov Decision Processes: Value Functions

We can now measure the *value* of a state/state-action pair under a policy $\pi \in \Pi$.

- ▶ Again, we appeal to the expected return...

State-Value

$$v_\pi^\gamma(x) \doteq \mathbb{E}_\pi[G_n^\gamma \mid X_n = x]$$

Action-Value

$$q_\pi^\gamma(x, u) \doteq \mathbb{E}_\pi[G_n^\gamma \mid X_n = x, U_n = u]$$

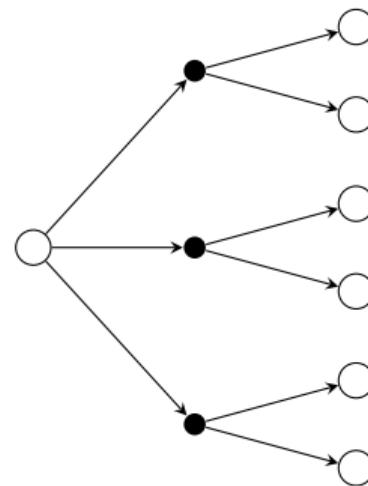
Markov Decision Processes: Bellman Equations

The value functions can be expressed recursively:

- ▶ This is the essence of *dynamic programming*.

Backup diagrams provide some insight:

- ▶ White nodes represent states.
- ▶ Black nodes represent actions.
- ▶ Edges represent probabilistic relations (transitions).

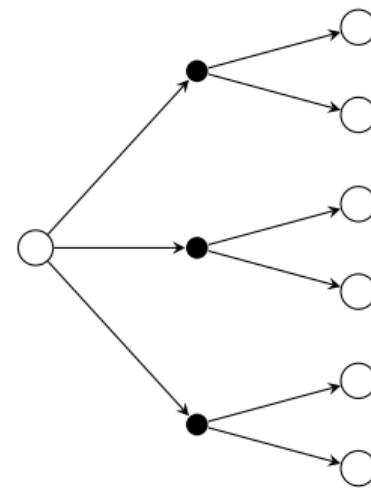


Markov Decision Processes: Bellman Equations

First define the *expected immediate reward*: $r(x, u) \doteq \mathbb{E}[R_{n+1} | X_n = x, U_n = u]$.

It then follows that:

$$\begin{aligned}
 q_\pi^\gamma(x, u) &= \mathbb{E}_\pi[G_n^\gamma | X_n = x, U_n = u], \\
 &= \mathbb{E}_\pi[R_{n+1} + \gamma G_{n+1}^\gamma | X_n = x, U_n = u], \\
 &= r(x, u) + \gamma \mathbb{E}_\pi[G_{n+1}^\gamma | X_n = x, U_n = u], \\
 &= \underbrace{r(x, u)}_{\text{Immediate}} + \underbrace{\gamma \mathbb{E}_\pi[v_\pi^\gamma(X_{n+1}) | X_n = x, U_n = u]}_{\text{Future}}.
 \end{aligned}$$



Markov Decision Processes: Optimality

Our objective can be expressed as: $J(\pi) = \mathbb{E}_\pi [v_\pi^\gamma(X_0) \mid X_0 \sim d_0]$.

Value functions thus define a *partial ordering over the set of policies* Π :

$$\pi \succeq \pi' \iff v_\pi^\gamma(x) \geq v_{\pi'}^\gamma(x) \quad \forall x \in \mathcal{X}.$$

Optimal Policy

$$\pi_* \in \arg \max_{\pi \in \Pi} J(\pi)$$

Optimal Value

$$v_*^\gamma(x) \doteq \max_{\pi \in \Pi} v_\pi^\gamma(x) = v_{\pi_*}^\gamma(x)$$

Markov Decision Processes: Optimality

What can we say about the set $\Pi_\star \doteq \arg \max_{\pi \in \Pi} J(\pi)$?

Existence

If values are bounded,¹ then $|\Pi_\star| > 0$.

Uniqueness

No guarantees in general!

¹Assured, for example, if $r : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \rightarrow \mathbb{R}$ is bounded and $\gamma \in [0, 1)$.

Markov Decision Processes: Optimality

The set Π_* also induces an optimal action-value function:

$$q_*^\gamma(x, u) \doteq \max_{\pi \in \Pi} q_\pi(x, u) = q_{\pi_*}(x, u) \quad \forall \pi_* \in \Pi_*$$

The functions v_*^γ and q_*^γ are then related by the equation,

$$v_*^\gamma(x) = \max_{u \in U} q_*^\gamma(x, u),$$

yielding the *Bellman Optimality equation*:

$$\begin{aligned} q_*^\gamma(x, u) &= r(x, u) + \gamma \mathbb{E}[v_*^\gamma(X_{n+1}) \mid X_n = x, U_n = u], \\ &= r(x, u) + \gamma \mathbb{E} \left[\max_{u' \in U} q_*^\gamma(X_{n+1}, u') \mid X_n = x, U_n = u \right]. \end{aligned}$$

Policy Evaluation

The policy evaluation problem² is stated as the following:

Given an MDP and a fixed policy $\pi \in \Pi$, compute the value functions $v_\pi^\gamma(x)$ and $q_\pi^\gamma(x, u)$ for all $x \in \mathbf{X}$ and $u \in \mathbf{U}$.

²Often referred to as the prediction problem.

Policy Evaluation

The “simplest” evaluation algorithm is **value iteration** (VI).³

- ▶ VI converges for fixed $\pi \in \Pi$:

$$\lim_{i \rightarrow \infty} \hat{v}_i \rightarrow v_\pi.$$

- ▶ Practically... it's not so good.
- ▶ Requires *full backups*.
- ▶ Uses *lots of memory*.

```

1 Initialise  $i \leftarrow 0, \hat{v}_i, \hat{q}_i$  arbitrarily;
2 while  $\|\hat{v}_i - v_\pi\|_\infty > \varepsilon$  do
3    $\hat{v}_{i+1} \leftarrow \hat{v}_i;$ 
4    $\hat{q}_{i+1} \leftarrow \hat{q}_i;$ 
5   for  $x \in \mathcal{X}$  do
6     for  $u \in \mathcal{U}$  do
7        $\hat{q}_{i+1}(x, u) \leftarrow r(x, u) +$ 
         $\gamma \sum_{x' \in \mathcal{X}} P(x, \{x'\} | u) \hat{v}_i(x');$ 
8     end
9      $\hat{v}_{i+1}(x) = \sum_{u \in \mathcal{U}} \pi(u | x) \hat{q}_{i+1}(x, u);$ 
10    end
11     $i \leftarrow i + 1;$ 
12 end

```

³Technically this isn't historically referred to as value iteration, but it's pretty close...

Policy Evaluation

What if we don't have direct access to P or r ?

Policy Evaluation

What if we don't have direct access to P or r ?

Sampling to the rescue!

Policy Evaluation

Temporal-difference (TD) learning uses *bootstrapping* to perform *sample backups*.

- ▶ TD converges for fixed $\pi \in \Pi$:

$$\lim_{i \rightarrow \infty} \hat{v}_i \rightarrow v_\pi.$$

- ▶ Often more practical than VI.
- ▶ *Variance is higher.*
- ▶ Sensitive to *sampling distribution*.

```

1 Sample  $x \sim d_0$ ;
2 Initialise  $i \leftarrow 0$ , and  $\hat{v}$  arbitrarily;
3 loop
4   Sample action:  $u \sim \pi(\cdot|x)$ ;
5   Transition to new state:  $x' \sim P(x, \cdot)$ ;
6   Compute the temporal-difference error:
7      $\delta \leftarrow r(x, u, x') + \gamma \hat{v}(x') - \hat{v}(x)$ ;
8   Update estimate:  $\hat{v}(x) \leftarrow \hat{v}(x) + \alpha \delta$ ;
9 end

```

Policy Evaluation: Approximate Methods

What now if \mathcal{U} is infinite? Or if no exact representation of v_π exists?

Policy Evaluation: Approximate Methods

What now if \mathcal{U} is infinite? Or if no exact representation of v_π exists?

(Stochastic) Gradient Descent to the rescue, as ever!

Policy Evaluation: Approximate Methods

Approximate methods resort to parameterised representations:

$$v_\pi(\cdot) \approx \hat{v}(\cdot; \mathbf{w}).$$

An important class is the set of *linear approximations*:

$$\mathcal{F}_\Phi = \left\{ f(\cdot) = \langle \mathbf{w}, \Phi(\cdot) \rangle = \sum_{i=1}^n w_i \phi_i(\cdot) \mid \mathbf{w} \in \mathbb{R}^n \right\},$$

for which

$$\nabla_{\mathbf{w}} f(\cdot; \mathbf{w}) = \Phi(\cdot) \quad \forall f \in \mathcal{F}_\Phi.$$

Policy Evaluation: Approximate Methods

There are lots of different choices for ϕ ... get creative!

Polynomial

$$\phi_i(\mathbf{x}) \doteq \prod_{j=1}^n x_j^{e_{ij}}$$

Fourier

$$\phi_i(\mathbf{x}) \doteq \cos(\pi \cdot \langle \mathbf{c}_i, \mathbf{x} \rangle)$$

Tile Coding

$$\phi_i(\mathbf{x}) \doteq \frac{1}{K} \sum_{j=1}^K \mathbf{1}(\text{Project}_j[\mathbf{x}] = i)$$

Policy Evaluation: Approximate Methods

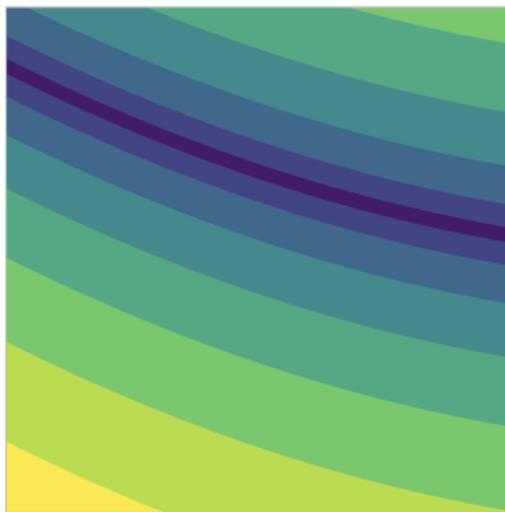


Figure: Ground truth.

\approx

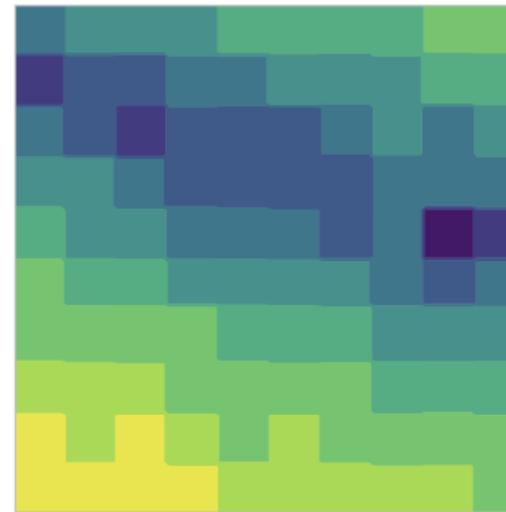


Figure: Tile Coding ($K = 1$)

Policy Evaluation: Approximate Methods

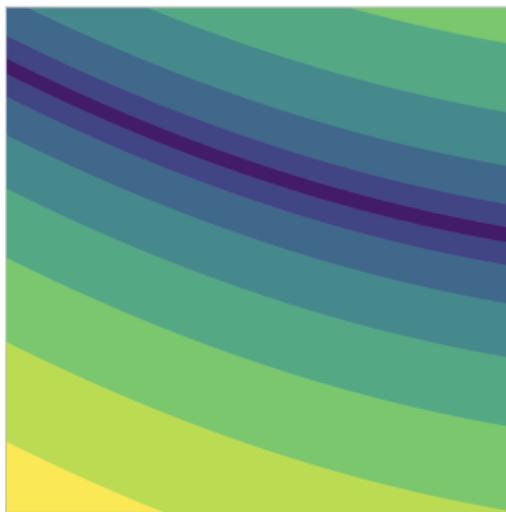
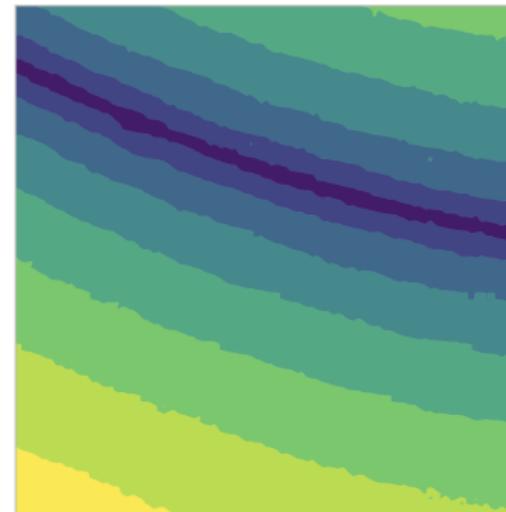
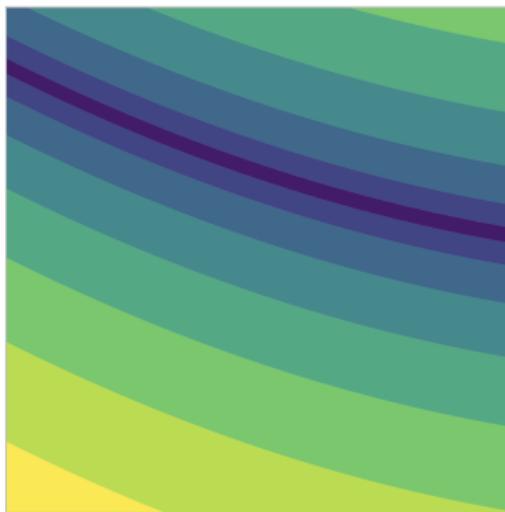
 \approx 

Figure: Ground truth.

Figure: Tile Coding ($K = 16$)

Policy Evaluation: Approximate Methods



≈

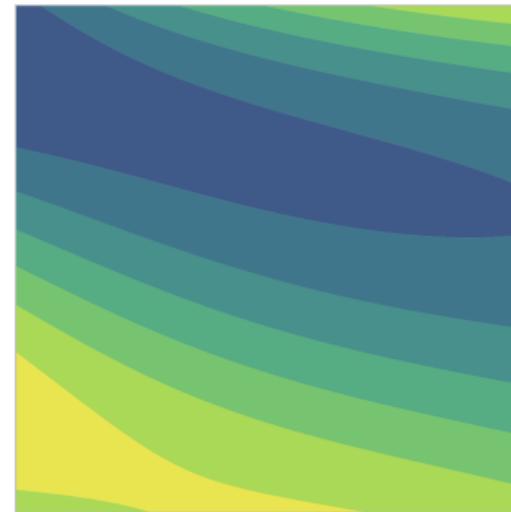


Figure: Ground truth.

Figure: Polynomial (3rd-order)

Policy Evaluation: Approximate Methods

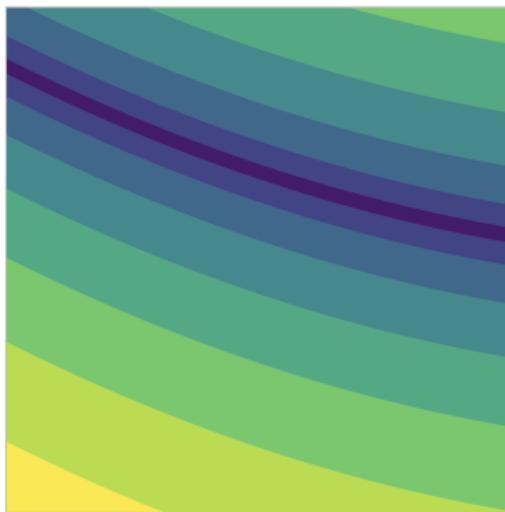


Figure: Ground truth.

\approx



Figure: Fourier (7th-order)

Policy Evaluation: Approximate Methods

Let us define an (*mean-squared*) *error* associated with the estimate \hat{v}_w ,

$$\mathcal{E}(\mathbf{w}) \doteq \|\hat{v}_w - v_\pi\|_d^2 = \mathbb{E}\left[[\hat{v}_w(X) - v_\pi(X)]^2 \mid X \sim d(\cdot) \right],$$

and then differentiate to give:

$$\nabla_{\mathbf{w}} \mathcal{E}(\mathbf{w}) = 2\mathbb{E}[[\hat{v}_w(X) - v_\pi(X)] \nabla_{\mathbf{w}} \hat{v}_w(X) \mid X \sim d(\cdot)].$$

Policy Evaluation: Approximate Methods

The corresponding *stochastic approximation*⁴ is given by

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \alpha_t [\hat{v}_{\mathbf{w}}(\mathbf{X}_n) - \underbrace{v_{\pi}(\mathbf{X}_n)}_{\text{Chicken } \iff \text{ Egg?}}] \nabla_{\mathbf{w}} \hat{v}_{\mathbf{w}}(\mathbf{X}).$$

Fortunately, we can replace $v_{\pi}(\mathbf{X}_n)$ with any *unbiased estimator* Z_n :

MC Let $Z_n \doteq \sum_{i=1}^{N-n} \gamma^{i-1} R_{n+i}$.

TD Let $Z_n \doteq R_{n+1} + \gamma \hat{v}(\mathbf{X}_{n+1})$.

⁴Vivek S Borkar and Sean P Meyn. "The ODE Method for Convergence of Stochastic Approximation and Reinforcement Learning". In: SIAM Journal on Control and Optimization 38.2 (2000), pp. 447–469.

Policy Evaluation

There have been many extensions to these methods:⁵

Least-Squares Methods⁶

Deep Learning

True-Online Methods⁷

Eligibility Traces⁸

Gradient Methods⁹

⁵ Adam White and Martha White. "Investigating Practical Linear Temporal Difference Learning". In: arXiv preprint arXiv:1602.08771 (2016).

⁶ Steven J Bradtke and Andrew G Barto. "Linear Least-Squares Algorithms for Temporal Difference Learning". In: Machine Learning 22.1 (1996), pp. 33–57

⁷ Harm Van Seijen et al. "True Online Temporal-Difference Learning". In: JMLR 17.1 (2016), pp. 5057–5096

⁸ Hamid Reza Maei. "Gradient Temporal-Difference Learning Algorithms". 2011

⁹ Richard S Sutton and Andrew G Barto. Reinforcement Learning: An Introduction. MIT press, 2018

Policy Optimisation

The policy optimisation problem¹⁰ is stated as the following:

Given an MDP, find an optimal policy $\pi_\star \in \Pi_\star$.

¹⁰Often referred to as the control problem.

Policy Optimisation: Value-Based Methods

The first class of algorithms are *value-based methods*.

The goal is to construct a sequence $(\hat{q}_0, \hat{q}_1, \hat{q}_2, \dots)$ such that

$$\lim_{n \rightarrow \infty} q_n = q_\star.$$

Since acting greedily wrt. q_\star is optimal, it follows that

$$\lim_{n \rightarrow \infty} \text{Greedy}[q_n] \in \Pi_\star.$$

Policy Optimisation: Value-Based Methods

We now use the VI algorithm from before, but for control.

- ▶ Note the use of the Bellman optimality equations!

```
1 Initialise  $i \leftarrow 0, \hat{v}_i(x), \hat{q}_i(x, u)$  arbitrarily;
2 while  $\|\hat{v}_i - v_\star\|_\infty > \varepsilon$  do
3    $\hat{v}_{i+1} \leftarrow \hat{v}_i;$ 
4    $\hat{q}_{i+1} \leftarrow \hat{q}_i;$ 
5   for  $x \in X$  do
6     for  $u \in U$  do
7        $\hat{q}_{i+1}(x, u) \leftarrow r(x, u) + \gamma \sum_{x' \in X} P(x, \{x'\} | u) \max_{u' \in U} \hat{q}_i(x', u');$ 
8     end
9      $\hat{v}_{i+1}(x) = \max_{u \in U} \hat{q}_{i+1}(x, u);$ 
10    end
11     $i \leftarrow i + 1;$ 
12 end
```

Policy Optimisation: Value-Based Methods

As with evaluation, we can *leverage sampling* when P or r are not available.

- ▶ The direct result is **Watkins' Q-Learning** algorithm.

```
1 Initialise  $x \sim d_0$ , and  $\hat{q}_w(x, u)$  arbitrarily;
2 loop
3   Sample action:  $u \sim \pi(\cdot|x)$ ;
4   Transition to new state:  $x' \sim P(x, \cdot)$ ;
5   Compute the temporal-difference error:  $\delta \leftarrow r(x, u, x') + \gamma \max_{u' \in U} \hat{q}(x', u') - \hat{q}(x, u)$ ;
6   Update weights:  $w \leftarrow w + \alpha \delta \nabla_w \hat{q}_w(x, u)$ ;
7   Innovate:  $x \leftarrow x'$ ;
8 end
```

Policy Optimisation: Policy-Gradient Methods

The second class of algorithms are *policy-gradient methods*.

Our goal here is to construct a sequence $(\theta_0, \theta_1, \theta_2, \dots)$ such that

$$\lim_{n \rightarrow \infty} \nabla_{\theta_n} J(\pi_{\theta_n}) = 0,$$

whereby, we hope,¹¹

$$\lim_{n \rightarrow \infty} \pi_{\theta_n} \in \Pi_\star.$$

¹¹This is not guaranteed!

Policy Optimisation: Policy-Gradient Methods

Let $d_\pi^\gamma \in \Delta(X)$ denote the (*improper*) *discounted occupancy measure*.

- ▶ This describes “how long we spend in each state (with discounting).”

The *policy gradient theorem*¹² then tells us that

$$\nabla_{\theta} J(\pi_{\theta}) = \underbrace{\mathbb{E}_{d_{\pi_{\theta}}^{\gamma}, \pi_{\theta}} [q_{\pi_{\theta}}(X, U) \nabla_{\theta} \log \pi_{\theta}(U | X)]}_{\text{Score-Based Gradient Estimator}},$$

which allows us to perform local search:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\pi_{\theta}).$$

¹²Richard S Sutton et al. “Policy Gradient Methods for Reinforcement Learning with Function Approximation”. In: Proc. of NeurIPS 12 (1999).

Policy Optimisation: Policy-Gradient Methods

There are lots of variants on policy gradient methods:

- ▶ Most revolve around how to reduce variance on estimates of $\nabla_{\theta} J(\pi_{\theta})$.

Baselines

Actor-Critic Methods

**Generalised Advantage
Estimation¹³**

SVRPG¹⁴

Topological¹⁵

¹³John Schulman et al. "High-Dimensional Continuous Control using Generalized Advantage Estimation". In: *arXiv preprint arXiv:1506.02438* (2015)

¹⁴Matteo Papini et al. "Stochastic Variance-Reduced Policy Gradient". In: *Proc. of ICML. 2018*, pp. 4026–4035

¹⁵Yasuhiro Fujita and Shin-ichi Maeda. "Clipped Action Policy Gradient". In: *Proc. of ICML. 2018*, pp. 1597–1606

Policy Optimisation: Policy-Gradient Methods

The *baseline trick* is particularly important (and will feature later):

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{d_{\pi_{\theta}}^{\gamma}, \pi_{\theta}} [[q_{\pi_{\theta}}(X, U) - b(X)] \nabla_{\theta} \log \pi_{\theta}(U | X)],$$

which holds because

$$\mathbb{E}_{d_{\pi_{\theta}}^{\gamma}, \pi_{\theta}} [b(X) \nabla_{\theta} \log \pi_{\theta}(U | X)] = 0.$$

The *optimal baseline is impractical* to derive or compute.

- ▶ Instead, we typically use $b(x) \doteq v_{\pi_{\theta}}(x)$.
- ▶ But there are better¹⁶ choices...

¹⁶If it wasn't already clear, this is a hint for later!

Algorithmic Trading

Trading Formalism

The trading universe comprises *one riskless asset (cash) and $m > 0$ risky assets.*

Trading Formalism

The trading universe comprises *one riskless asset (cash) and $m > 0$ risky assets.*

Riskless Asset

- ▶ Price fixed at unity.
- ▶ Agent holdings are denoted by

$$\mathbf{C} \doteq (C_0, C_1, C_2, \dots).$$

Trading Formalism

The trading universe comprises *one riskless asset (cash) and $m > 0$ risky assets.*

Riskless Asset

- ▶ Price fixed at unity.
- ▶ Agent holdings are denoted by

$$\mathbf{C} \doteq (C_0, C_1, C_2, \dots).$$

Risky Assets

- ▶ Prices \mathbf{Z} evolve stochastically with
$$\mathbf{Z}_n = [(\mathbf{Z}_n)_1, (\mathbf{Z}_n)_2, \dots, (\mathbf{Z}_n)_m] \in \mathbb{R}^m.$$
- ▶ Agent holdings are Ω with
$$\Omega_n = [(\Omega_n)_1, (\Omega_n)_2, \dots, (\Omega_n)_m] \in \mathbb{R}^m.$$

Trading Formalism

The agent portfolio is denoted by the tuple $\Lambda_n \doteq (C_n, \Omega_n)$ for all $n \geq 0$.

The *mark-to-market value* of Λ_n is the inner product:

$$\Upsilon_n \doteq C_n + \langle \Omega_n, Z_n \rangle .$$

Trading Formalism

The agent portfolio is denoted by the tuple $\Lambda_n \doteq (C_n, \Omega_n)$ for all $n \geq 0$.

The *mark-to-market value* of Λ_n is the inner product:

$$\Upsilon_n \doteq C_n + \langle \Omega_n, Z_n \rangle .$$

We can then compute the discrete derivative as follows:

$$\begin{aligned} \Delta \Upsilon_n &= \Delta C_n + \Delta \langle \Omega_n, Z_n \rangle , \\ &= \Delta C_n + \langle \Delta \Omega_n, Z_n \rangle + \langle \Delta \Omega_n, \Delta Z_n \rangle + \langle \Omega_n, \Delta Z_n \rangle , \\ &= \underbrace{\Delta C_n + \langle \Delta \Omega_n, Z_n \rangle}_{\text{Trading}} + \underbrace{\langle \Omega_{n+1}, \Delta Z_n \rangle}_{\text{Speculation}} . \end{aligned}$$

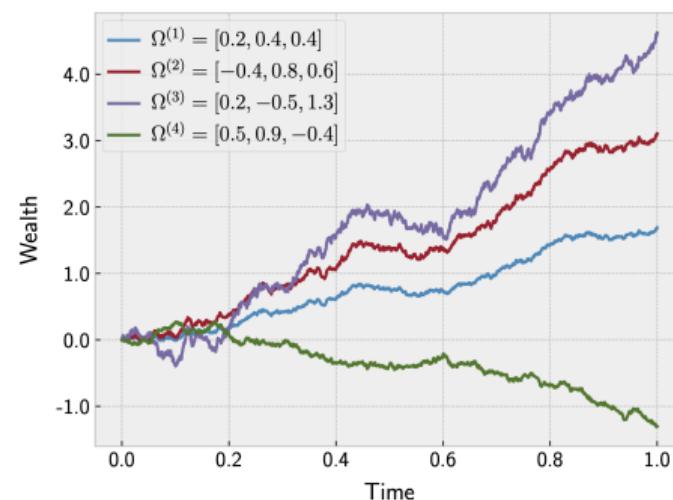
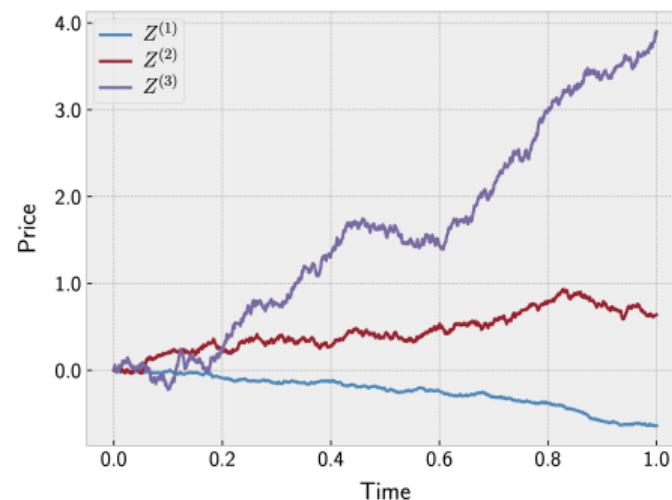
Note the $n + 1!$

As an example, consider a portfolio consisting of three risky assets:

$$\Delta \mathbf{Z}_n = \begin{bmatrix} -0.5 \\ 0.5(1.5 - Z_n^{(2)}) \\ 3 \end{bmatrix} \Delta t + \begin{bmatrix} 5 \\ 15 \\ 25 \end{bmatrix} \circ \mathbf{N}_n.$$

As an example, consider a portfolio consisting of three risky assets:

$$\Delta Z_n = \begin{bmatrix} -0.5 \\ 0.5(1.5 - Z_n^{(2)}) \\ 3 \end{bmatrix} \Delta t + \begin{bmatrix} 5 \\ 15 \\ 25 \end{bmatrix} \circ N_n.$$



Trading Formalism: Transactions

Of course, we can do more than just buy-and-hold!

- ▶ So lets enrich our model to include transactions:

Sell

Let $\nu^- \geq 0$ denote the agent's
ask transaction flow.

Buy

Let $\nu^+ \geq 0$ denote the agent's
bid transaction flow.

Trading Formalism: Transactions

Of course, we can do more than just buy-and-hold!

- ▶ So lets enrich our model to include transactions:

Sell

Let $\mathbf{v}^- \geq 0$ denote the agent's
ask transaction flow.

Buy

Let $\mathbf{v}^+ \geq 0$ denote the agent's
bid transaction flow.

We can now express the inventory process:

$$\Delta \boldsymbol{\Omega}_n = \mathbf{v}_n^+ - \mathbf{v}_n^- \quad \Rightarrow \quad \boldsymbol{\Omega}_n = \omega_0 + \sum_{k=0}^{n-1} \mathbf{v}_k^+ - \mathbf{v}_k^-$$

Trading Formalism: Transactions

Of course, we can do more than just buy-and-hold!

- So lets enrich our model to include transactions:

Sell

Let $\mathbf{v}^- \geq 0$ denote the agent's
ask transaction flow.

Buy

Let $\mathbf{v}^+ \geq 0$ denote the agent's
bid transaction flow.

And its positive/negative parts:

$$\Delta\Omega_n^\pm = \mathbf{v}_n^\pm \quad \Rightarrow \quad \Omega_n^\pm = \omega_0^\pm + \sum_{k=0}^{n-1} \mathbf{v}_k^\pm \quad \Rightarrow \quad \Omega_n = \Omega_n^+ - \Omega_n^-$$

Trading Formalism: Transactions

So, what happens when the agent *buys k units of a risky asset*?

Agent

- ▶ Sends kZ_n cash to counterparty.
 - ▶ $C_{n+1} \leftarrow C_n - kZ_n.$
 - ▶ $\Omega_{n+1} \leftarrow \Omega_n + k.$
 - ▶ $v_n = v_n^+ = k.$

Counterparty

- ▶ Sends k units to agent.

Trading Formalism: Transactions

So, what happens when the agent *sells k units of a risky asset*?

Agent

- ▶ Sends k units to counterparty.
 - ▶ $C_{n+1} \leftarrow C_n + kZ_n.$
 - ▶ $\Omega_{n+1} \leftarrow \Omega_n - k.$
 - ▶ $v_n = -v_n^- = -k.$

Counterparty

- ▶ Sends kZ_n cash to agent.

Trading Formalism: Transactions

Is it realistic to assume that you can trade at Z_n ?

Trading Formalism: Transactions

Is it realistic to assume that you can trade at Z_n ? No!

- ▶ Latency.
- ▶ Slippage.
- ▶ Transaction fees.
- ▶ Information asymmetry.

Trading Formalism: Transactions

Is it realistic to assume that you can trade at Z_n ? No!

- ▶ Latency.
- ▶ Slippage.
- ▶ Transaction fees.
- ▶ Information asymmetry.

We can account for some effects with a (*temporary*) *impact function* $f : \dots \rightarrow \mathbb{R}$

$$\Delta C_n \approx -Z_n \cdot v_n + f(Z_n, v_n^+, v_n^-, \dots),$$

but reality is a little more complicated...

Limit Order Books

Many electronic markets operate as *limit order books (LOB)*.

A *limit order (LO)* is a triple:

$$o \doteq (n_o, \omega_o, z_o).$$

The LOB is then the collection

$$\mathcal{B}_n \doteq \mathcal{B}_n^+ \cup \mathcal{B}_n^-$$

where

$$\mathcal{B}_n^+ \doteq \{o : n_o \leq n, \omega_o > 0\},$$

$$\mathcal{B}_n^- \doteq \{o : n_o \leq n, \omega_o < 0\}.$$

Price	Offers
101.0	12
100.5	13
100.0	
35	99.5
11	99.0
3	98.5
Bids	

Limit Order Books

Many electronic markets operate as *limit order books (LOB)*.

A market order is a special case.

- ▶ Let the price tend to 0 or ∞ .

The process is called: *walking the book*.

- ▶ Transact until liquidity is exhausted.
- ▶ Or until the request is satisfied.

Price	
Offers	
Bids	
101.0	12
100.5	13
100.0	
99.5	35
99.0	11
98.5	3

Limit Order Books

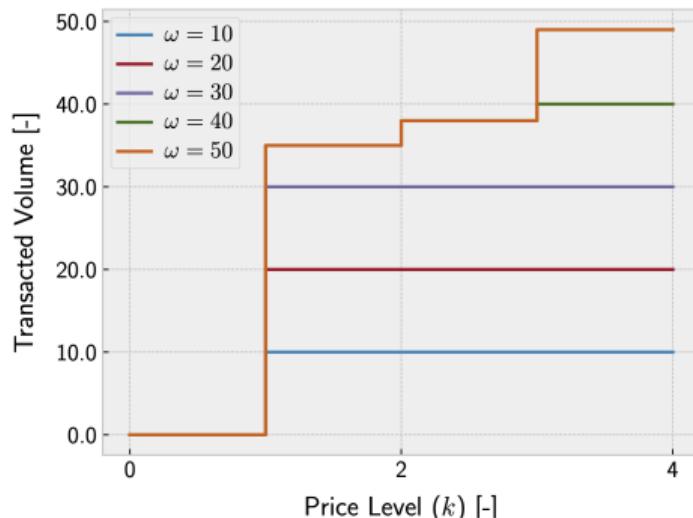
Many electronic markets operate as *limit order books (LOB)*.

A market order (MO) is a special case.

- ▶ Let the price tend to 0 or ∞ .

The process is called: *walking the book*.

- ▶ Transact until liquidity is exhausted.
- ▶ Or until the request is satisfied.



Limit Order Books

Why do you favour LOs over MOs?

Limit Order

- ▶ Passive execution.
- ▶ Better price.

Market Order

- ▶ Immediate execution.
- ▶ Incurs slippage.

Market Making

Market making is a canonical problem in algorithmic trading.

- ▶ Market makers (MMs) quote bid and ask prices.
- ▶ They provide liquidity to a market.
- ▶ Derive profit of transactions.

Market Making

Market making is a canonical problem in algorithmic trading.

- ▶ Market makers (MMs) quote bid and ask prices.
- ▶ They provide liquidity to a market.
- ▶ Derive profit of transactions.

There are some significant challenges for MMs:

- ▶ Balancing transaction rate and spread.
 - ▶ *Optimal Pricing*
- ▶ Avoiding exposure to large market movements.
 - ▶ *Inventory Risk*

Market Making: Avellaneda-Stoikov Model

The Avellaneda-Stoikov model is a parametric approximation of an LOB.

Agent trades in *unit LOs with prices about the mid-price*: $Z_n^\pm = Z_n \mp \delta_n^\pm$.

Market Making: Avellaneda-Stoikov Model

The Avellaneda-Stoikov model is a parametric approximation of an LOB.

Agent trades in *unit LOs with prices about the mid-price*: $Z_n^\pm = Z_n \mp \delta_n^\pm$.

Prices

Evolve as a random walk.

$$Z_0 = z_0,$$

$$Z_n = Z_{n-1} + \sigma N_n.$$

Market Making: Avellaneda-Stoikov Model

The Avellaneda-Stoikov model is a parametric approximation of an LOB.

Agent trades in *unit LOs with prices about the mid-price*: $Z_n^\pm = Z_n \mp \delta_n^\pm$.

Prices

Evolve as a random walk.

$$Z_0 = z_0,$$

$$Z_n = Z_{n-1} + \sigma N_n.$$

Transactions

Modelled as Poisson processes:

$$\lambda^\pm(\delta^\pm) = A^\pm \exp(-k^\pm \delta^\pm).$$

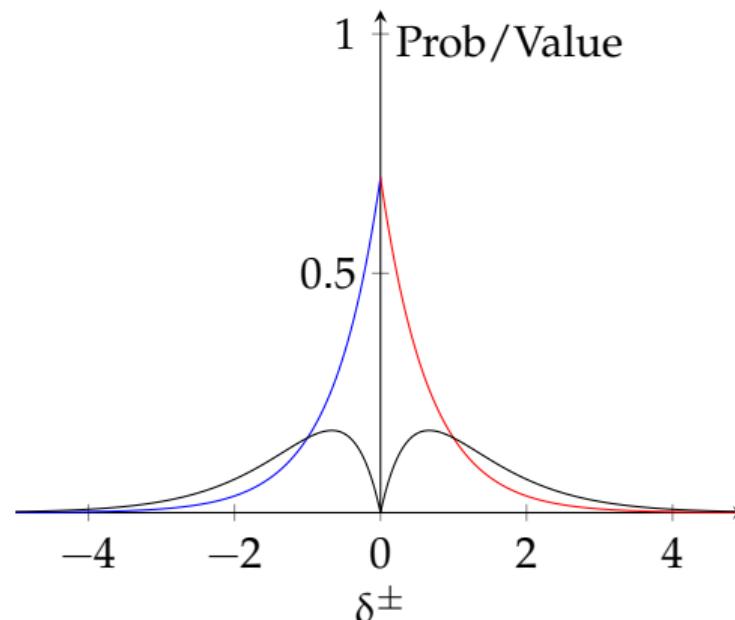
Market Making: Avellaneda-Stoikov Model

The probability of transacting decreases as δ^+ , δ^- grow.

The *expected value of transacting is (quasi) concave!*

- ▶ There is an “optimal” price.

But we also have to worry about inventory...



Market Making: (Quadratic) MDP Formulation

We begin by defining states/actions:

$$X_n \doteq (n, Z_n, \Omega_n^+, \Omega_n^-) \quad \text{and} \quad U_n \doteq (\delta_n^+, \delta_n^-).$$

Market Making: (Quadratic) MDP Formulation

We begin by defining states/actions:

$$X_n \doteq (n, Z_n, \Omega_n^+, \Omega_n^-) \quad \text{and} \quad U_n \doteq (\delta_n^+, \delta_n^-).$$

Letting $x_0 \doteq (0, z_0, 0, 0)$, we define the objective as

$$\begin{aligned} J(\pi) &= \mathbb{E}_\pi \left[\gamma_N - \beta \sum_{n=0}^{N-1} \Omega_{n+1}^2 \Delta Z_n \middle| X_0 = x_0 \right], \\ &= \mathbb{E}_\pi \left[\sum_{n=0}^{N-1} \Delta C_n + Z_n \Delta \Omega_n + \Omega_{n+1} \Delta Z_n (1 - \beta \Omega_{n+1}) \middle| X_0 = x_0 \right]. \end{aligned}$$

Market Making: (Quadratic) MDP Formulation

Observe that our cash process must evolve according to

$$\begin{aligned}\Delta C_n &= (Z_n + \delta_n^-) \Delta \Omega_n^- - (Z_n - \delta_n^+) \Delta \Omega_n^+, \\ &= \delta_n^- \Delta \Omega_n^- + \delta_n^+ \Delta \Omega_n^+ - Z_n \Delta \Omega_n.\end{aligned}$$

Market Making: (Quadratic) MDP Formulation

Observe that our cash process must evolve according to

$$\begin{aligned}\Delta C_n &= (Z_n + \delta_n^-) \Delta \Omega_n^- - (Z_n - \delta_n^+) \Delta \Omega_n^+, \\ &= \delta_n^- \Delta \Omega_n^- + \delta_n^+ \Delta \Omega_n^+ - Z_n \Delta \Omega_n.\end{aligned}$$

We can thus simplify the objective to give:

$$J(\pi) = \mathbb{E}_\pi \left[\sum_{n=0}^{N-1} \delta_n^- \Delta \Omega_n^- + \delta_n^+ \Delta \Omega_n^+ + \textcolor{red}{\Omega_{n+1} \Delta Z_n (1 - \beta \Omega_{n+1})} \middle| X_0 = x_0 \right].$$

Market Making: (Quadratic) MDP Formulation

One possible instantiation of the reward is thus

$$R_{n+1} = \delta_n^- \Delta \Omega_n^- + \delta_n^+ \Delta \Omega_n^+ + \Omega_{n+1} \Delta Z_n (1 - \beta \Omega_{n+1}).$$

From this it follows that

$$J(\pi) = \mathbb{E}_\pi \underbrace{\left[\sum_{n=0}^{N-1} R_{n+1} \mid X_0 = x_0 \right]}_{\text{Look familiar?!"}}$$

Market Making: (Quadratic) MDP Formulation

One possible instantiation of the reward is thus

$$R_{n+1} = \delta_n^- \Delta \Omega_n^- + \delta_n^+ \Delta \Omega_n^+ + \Omega_{n+1} \Delta Z_n (1 - \beta \Omega_{n+1}).$$

From this it follows that

$$J(\pi) = \mathbb{E}_\pi \underbrace{\left[\sum_{n=0}^{N-1} R_{n+1} \mid X_0 = x_0 \right]}_{\text{Look familiar?}} = v_\pi^1(x_0).$$

Market Making: Heuristic Strategies

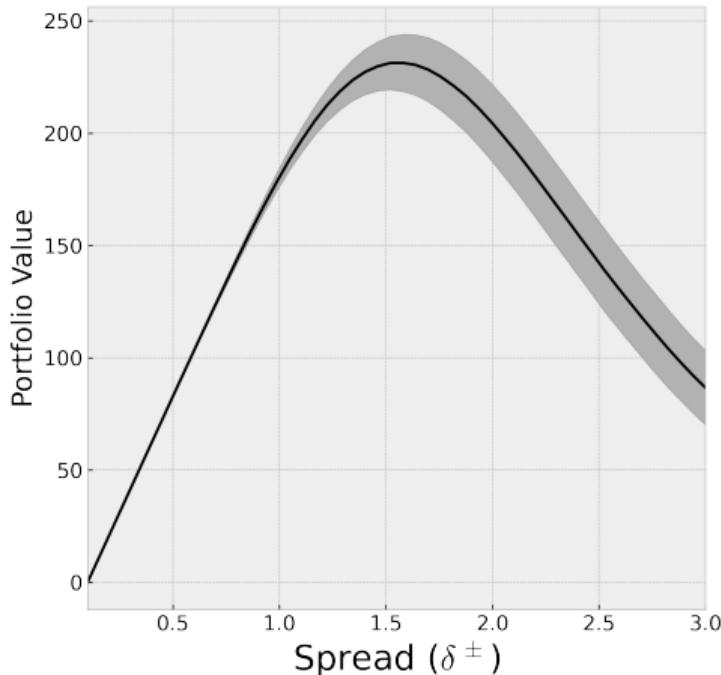
What if we always quote a fixed spread?

$$\delta_n^\pm \doteq \delta^\pm, \quad \forall n \geq 0.$$

Market Making: Heuristic Strategies

What if we always quote a fixed spread?

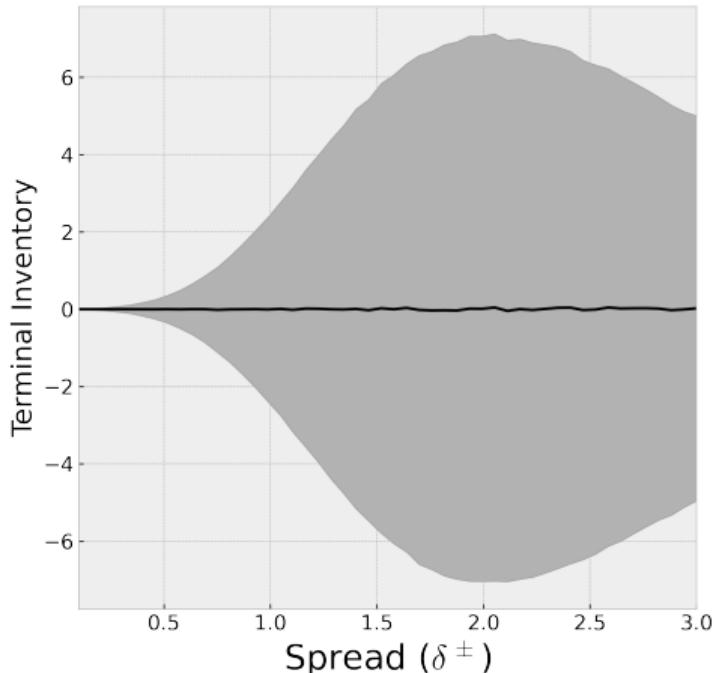
$$\delta_n^\pm \doteq \delta^\pm, \quad \forall n \geq 0.$$



Market Making: Heuristic Strategies

What if we always quote a fixed spread?

$$\delta_n^{\pm} \doteq \delta^{\pm}, \quad \forall n \geq 0.$$



Market Making: Heuristic Strategies

Can we avoid large positions?

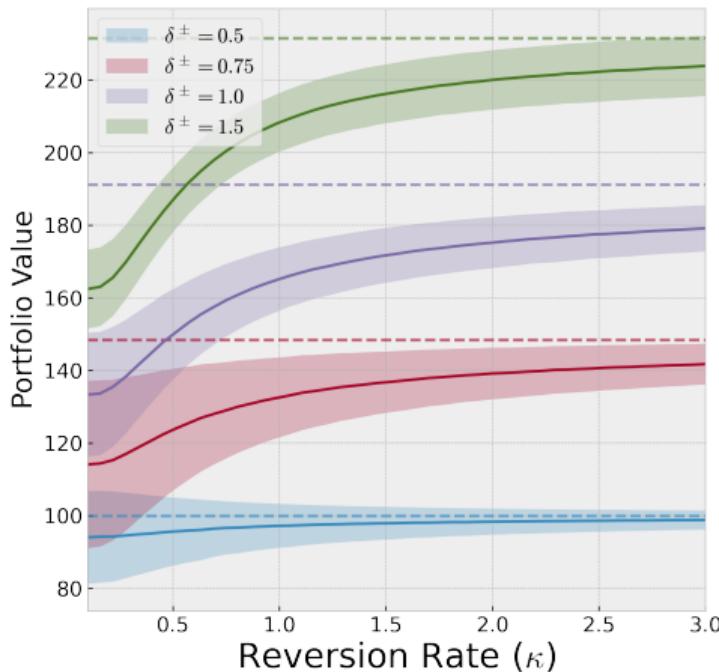
$$\delta_n^\pm \doteq \delta^\pm \cdot \exp\left(-\frac{1}{\kappa} [\Omega_n]_\pm\right).$$

Market Making: Heuristic Strategies

Can we avoid large positions?

$$\delta_n^\pm \doteq \delta^\pm \cdot \exp\left(-\frac{1}{\kappa} [\Omega_n]_\pm\right).$$

\implies Induce mean reversion in Ω_n .

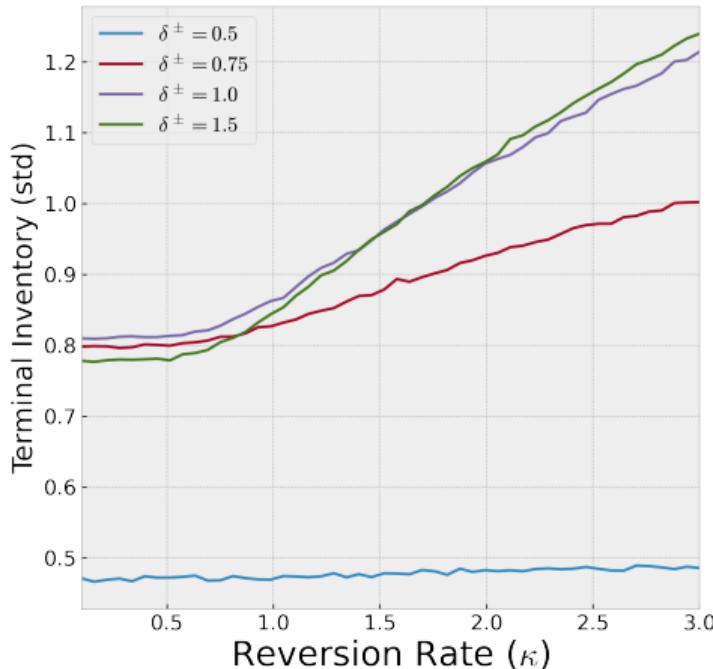


Market Making: Heuristic Strategies

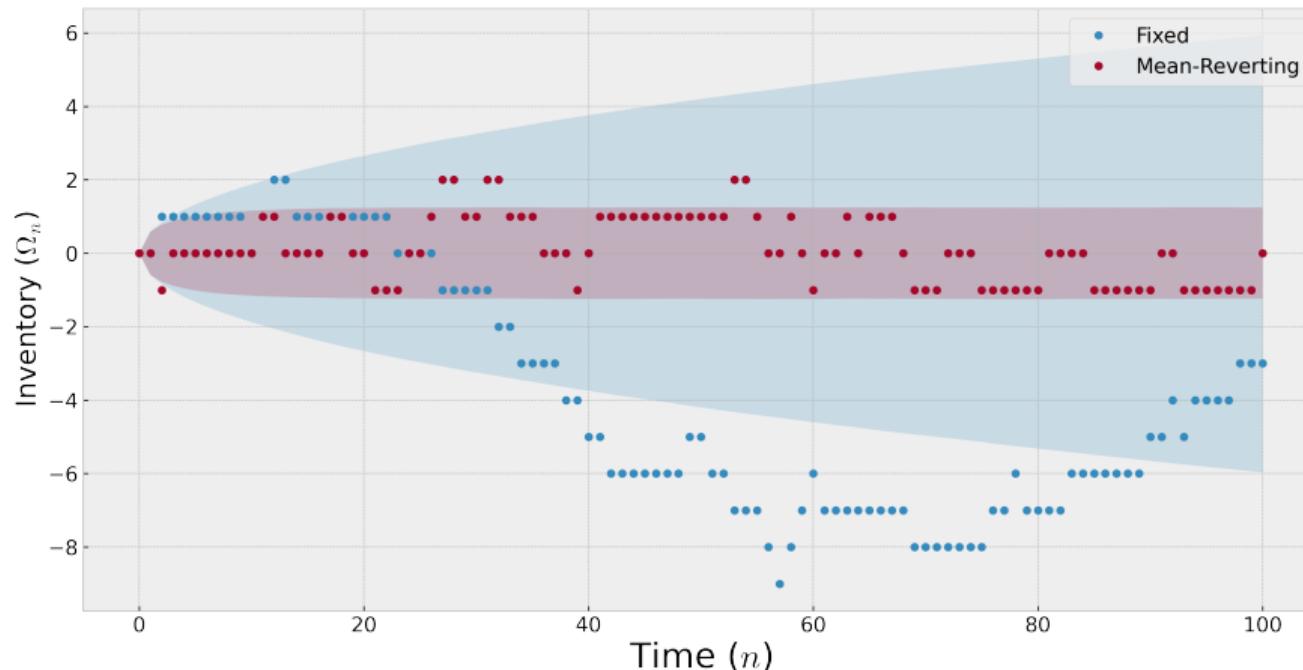
Can we avoid large positions?

$$\delta_n^\pm \doteq \delta^\pm \cdot \exp\left(-\frac{1}{\kappa} [\Omega_n]_\pm\right).$$

\implies Induce mean reversion in Ω_n .



Market Making: Heuristic Strategies



Market Making: Learnt Strategies

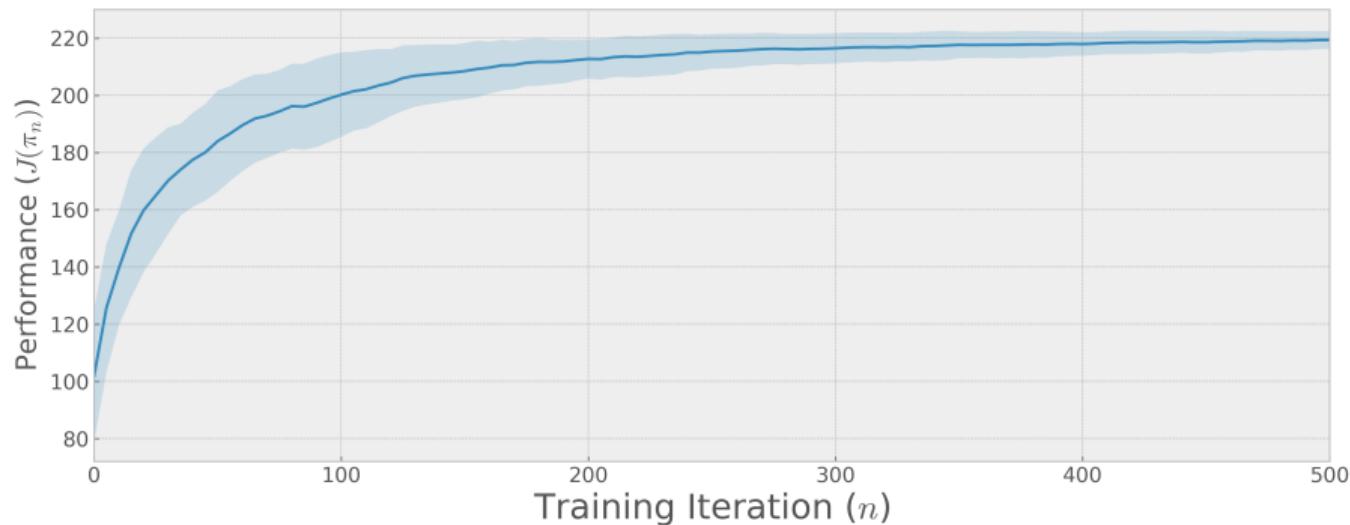
Alternatively, we can learn a policy!¹⁷

- ▶ Represent π_θ using a cubic polynomial.
 - ▶ Yes, this means we only need 8 parameters!
- ▶ Use the REINFORCE algorithm:
 - ▶ Generate 200 sample paths.
 - ▶ Take the average gradient for better stability.

¹⁷Happy to provide the code :)

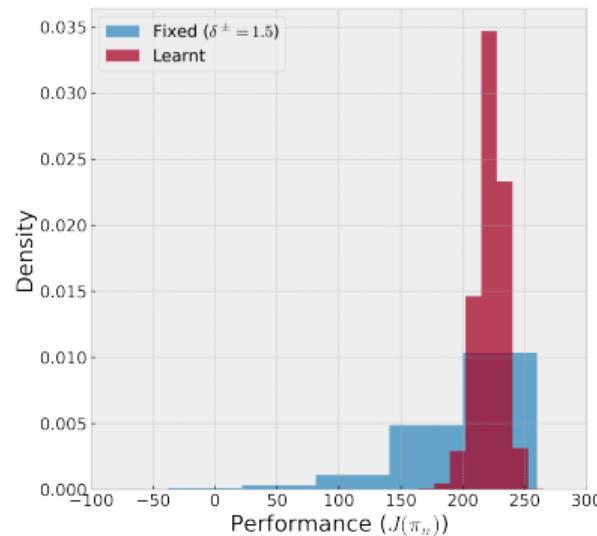
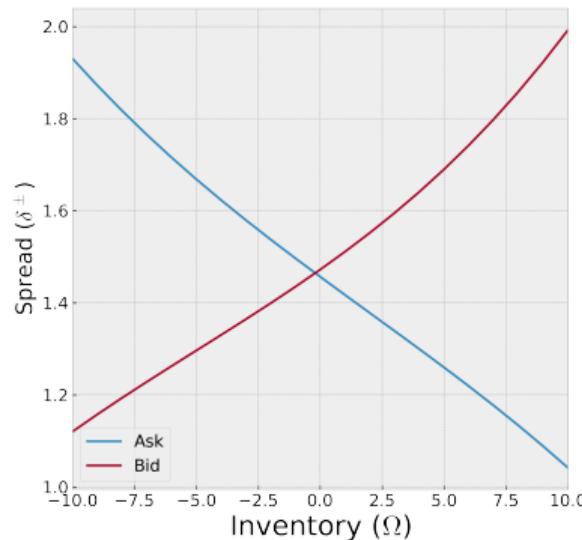
Market Making: Learnt Strategies

The performance improves smoothly over time (std over random seeds)!



Market Making: Learnt Strategies

The policy is smooth and interpretable, and also outperforms a fixed benchmark!



Gradient Factorisation

Scalability Issues

Policy gradient methods do not scale.

- ▶ We have to leverage prior knowledge!

Action-Dimensionality

If the dimensionality of U is large,
meaningful exploration is hard.

Objective-Multiplicity

If $J(\pi)$ is multi-objective, you get
high variance and chattering.

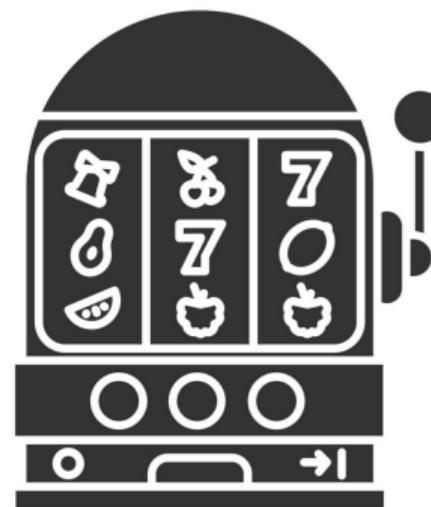
Scalability Issues

Let's say you're faced with n slot machines...

- ▶ you select $m \leq n$ arms...
- ▶ and see a **scalar reward** of $R = 5$.

Which arms were responsible?

- ▶ What if all arms gave positive reward?
- ▶ What if only one arm gave positive reward?



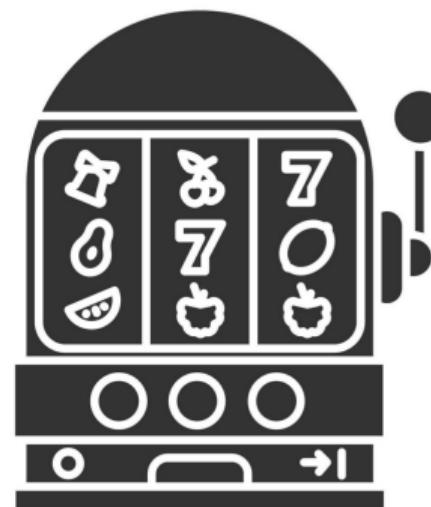
Scalability Issues

In reality, we see a vector, say:

$$\mathbf{R} = [0, \dots, 0, 1, 0, \dots 0].$$

The standard MDP formulation ignores this!

- ▶ So how do we leverage this structure?



Scalability Issues

Market Making

Cross-Currency Arbitrage

**Multi-Venue Optimal
Execution**

Water Reservoir Control

Elevator/Job Scheduling

**Energy Consumption
Optimisation**

Influence Networks

Consider a *scalarised multi-objective MDP* where

$$\begin{aligned} J(\pi) &= \mathbb{E}_\pi \left[q_\pi(X_0, U_0) \doteq \langle \boldsymbol{\lambda}, \mathbf{q}_\pi(X_0, U_0) \rangle = \sum_{i=1}^m \lambda_i q_{\pi,i}(X_0, U_0) \middle| X_0 \sim d_0 \right], \\ &= \sum_{i=1}^m \lambda_i \mathbb{E}_\pi [q_{\pi,i}(X_0, U_0) \mid X_0 \sim d_0]. \end{aligned}$$

Influence Networks

Consider a *scalarised multi-objective MDP* where

$$\begin{aligned} J(\pi) &= \mathbb{E}_\pi \left[q_\pi(X_0, U_0) \doteq \langle \boldsymbol{\lambda}, \mathbf{q}_\pi(X_0, U_0) \rangle = \sum_{i=1}^m \lambda_i q_{\pi,i}(X_0, U_0) \middle| X_0 \sim d_0 \right], \\ &= \sum_{i=1}^m \lambda_i \mathbb{E}_\pi [q_{\pi,i}(X_0, U_0) \mid X_0 \sim d_0]. \end{aligned}$$

The value function **breaks down into m distinct terms**:

- ▶ This will be familiar to those who've studied *value decomposition networks*.

Influence Networks

Consider a *scalarised multi-objective MDP* where

$$\begin{aligned}
 J(\pi) &= \mathbb{E}_\pi \left[q_\pi(X_0, U_0) \doteq \langle \boldsymbol{\lambda}, \mathbf{q}_\pi(X_0, U_0) \rangle = \sum_{i=1}^m \lambda_i q_{\pi,i}(X_0, \sigma_i(U_0)) \mid X_0 \sim d_0 \right], \\
 &= \sum_{i=1}^m \lambda_i \mathbb{E}_\pi [q_{\pi,i}(X_0, \sigma_i(U_0)) \mid X_0 \sim d_0].
 \end{aligned}$$

What if each term $q_{\pi,i}$ only depends on a *subset of the action U* ?

Influence Networks

This suggests a *linear-algebraic reformulation*:

$$J(\pi) = \mathbb{E}_\pi[q_\pi(X_0, U_0) \doteq K\lambda \circ q_\pi(X_0, U_0) \mid X_0 \sim d_0],$$

where

$$\begin{aligned} K &\in \mathbb{B}^{m \times m} & \text{and} & \lambda \doteq [\lambda_1, \dots, \lambda_m] \in \mathbb{R}^m, \\ q_\pi(x, u) &\doteq [q_{\pi,1}(x, u), \dots, q_{\pi,m}(x, u)] \in \mathbb{R}^m. \end{aligned}$$

Influence Networks

This suggests a *linear-algebraic reformulation*:

$$J(\pi) = \mathbb{E}_\pi[q_\pi(X_0, U_0) \doteq \mathbf{K}\boldsymbol{\lambda} \circ \mathbf{q}_\pi(X_0, U_0) \mid X_0 \sim d_0],$$

where

$$\underbrace{\mathbf{K} \in \mathbb{B}^{m \times m}}_{\text{Influence Matrix}} \quad \text{and} \quad \boldsymbol{\lambda} \doteq [\lambda_1, \dots, \lambda_m] \in \mathbb{R}^m,$$

$$\mathbf{q}_\pi(x, u) \doteq [q_{\pi,1}(x, u), \dots, q_{\pi,m}(x, u)] \in \mathbb{R}^m.$$

Influence Networks

Example

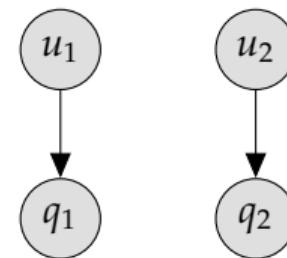
Consider a two-dimensional bandit of the form:

$$q(\cdot, \mathbf{u}) \doteq q_1(\cdot, u_1) + q_2(\cdot, u_2) = -|u_1 - c_1| - |u_2 - c_2|.$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Influence Network



Influence Networks

Example

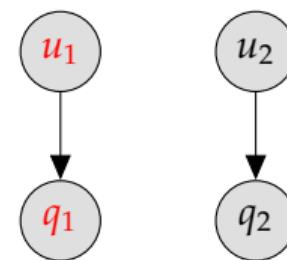
Consider a two-dimensional bandit of the form:

$$q(\cdot, \mathbf{u}) \doteq q_1(\cdot, u_1) + q_2(\cdot, u_2) = -|u_1 - c_1| - |u_2 - c_2|.$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Influence Network



Influence Networks

Example

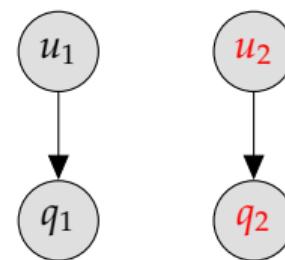
Consider a two-dimensional bandit of the form:

$$q(\cdot, \mathbf{u}) \doteq q_1(\cdot, u_1) + q_2(\cdot, u_2) = -|u_1 - c_1| - |u_2 - c_2|.$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Influence Network



Influence Networks

Example

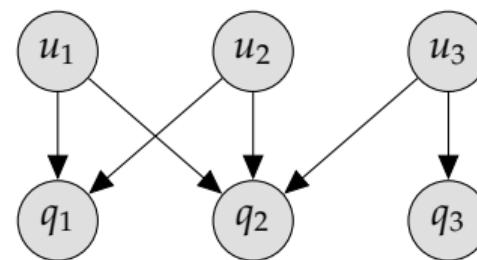
Consider an MDP with cross-dependency structure:

$$q(x, \mathbf{u}) \doteq \lambda_1 q_1(x, (u_1, u_2)) + \lambda_2 q_2(x, \mathbf{u}) + \lambda_3 q_3(x, u_3).$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

Influence Network



Influence Networks

Example

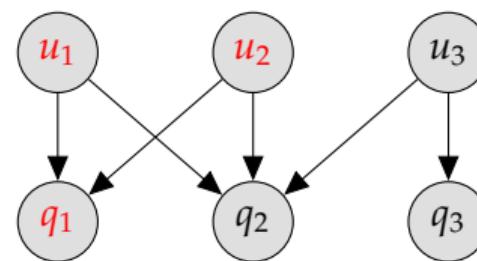
Consider an MDP with cross-dependency structure:

$$q(x, \mathbf{u}) \doteq \lambda_1 q_1(x, (u_1, u_2)) + \lambda_2 q_2(x, \mathbf{u}) + \lambda_3 q_3(x, u_3).$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

Influence Network



Influence Networks

Example

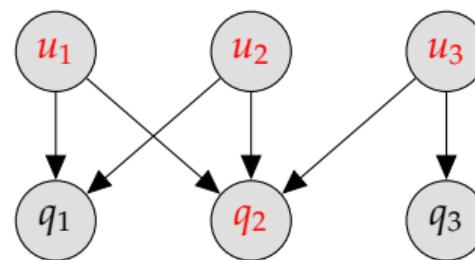
Consider an MDP with cross-dependency structure:

$$q(x, \mathbf{u}) \doteq \lambda_1 q_1(x, (u_1, u_2)) + \lambda_2 q_2(x, \mathbf{u}) + \lambda_3 q_3(x, u_3).$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

Influence Network



Influence Networks

Example

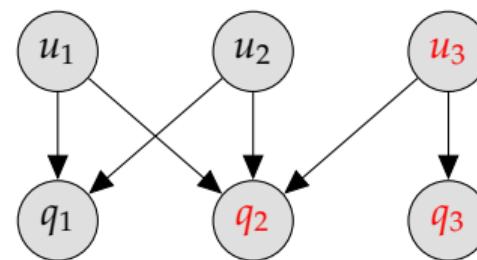
Consider an MDP with cross-dependency structure:

$$q(x, \mathbf{u}) \doteq \lambda_1 q_1(x, (u_1, u_2)) + \lambda_2 q_2(x, \mathbf{u}) + \lambda_3 q_3(x, u_3).$$

Influence Matrix

$$\mathbf{K} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

Influence Network



Policy Factorisation

So how do we embed this knowledge in our algorithms?

¹⁸*For an MDP, we can always do this without affecting the optimum, though we may lose some generality...*

Policy Factorisation

So how do we embed this knowledge in our algorithms?

- ▶ We don't train actions, we train policies...

¹⁸For an MDP, we can always do this without affecting the optimum, though we may lose some generality...

Policy Factorisation

So how do we embed this knowledge in our algorithms?

- ▶ We don't train actions, we train policies...
- ▶ So we must first *impose the independence in our policy distribution.*

¹⁸For an MDP, we can always do this without affecting the optimum, though we may lose some generality...

Policy Factorisation

So how do we embed this knowledge in our algorithms?

- ▶ We don't train actions, we train policies...
- ▶ So we must first *impose the independence in our policy distribution.*

Using the σ /subset notation, we express the policy as a product distribution:¹⁸

$$\pi(\mathbf{u} \mid x; \boldsymbol{\theta}) = \prod_{i=1}^N \pi_i(\sigma_i^\pi(\mathbf{u}) \mid x; \boldsymbol{\theta}).$$

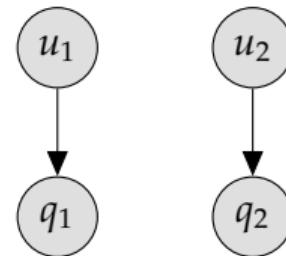
¹⁸For an MDP, we can always do this without affecting the optimum, though we may lose some generality...

Policy Factorisation

Recall the two-dimensional bandit:

$$q(\cdot, \mathbf{u}) = q_1(\cdot, u_1) + q_2(\cdot, u_2).$$

Influence Network



Policy Factorisation

Recall the two-dimensional bandit:

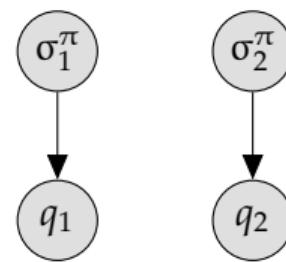
$$q(\cdot, \mathbf{u}) = q_1(\cdot, u_1) + q_2(\cdot, u_2).$$

We can naturally separate into:

$$\sigma_1^\pi(\mathbf{u}) = (u_1),$$

$$\sigma_2^\pi(\mathbf{u}) = (u_2).$$

Factored Influence Network

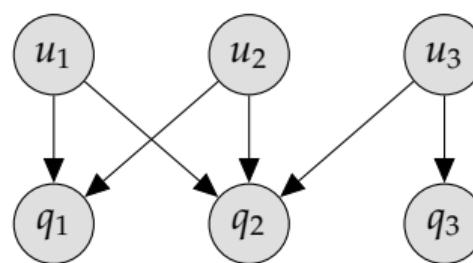


Policy Factorisation

Recall the mixed case:

$$\lambda_1 q_1(x, (u_1, u_2)) + \lambda_2 q_2(x, \mathbf{u}) + \lambda_3 q_3(x, u_3).$$

Influence Network



Policy Factorisation

Recall the mixed case:

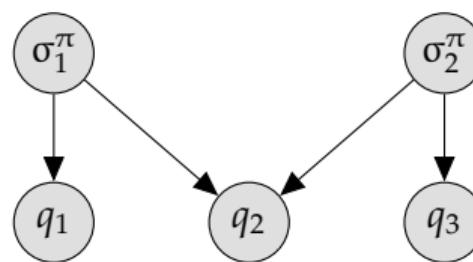
$$\lambda_1 q_1(x, (u_1, u_2)) + \lambda_2 q_2(x, \mathbf{u}) + \lambda_3 q_3(x, u_3).$$

We can naturally separate into:

$$\sigma_1^\pi(\mathbf{u}) = (u_1, u_2),$$

$$\sigma_2^\pi(\mathbf{u}) = (u_3).$$

Factored Influence Network



Gradient Factorisation

Recall Sutton's vanilla policy gradient:¹⁹

$$\begin{aligned}\nabla_{\theta} J(\pi) &= \mathbb{E}_{\pi}[\nabla_{\theta} \log \pi(\mathbf{U} | X) q(X, \mathbf{U})], \\ &\doteq \mathbb{E}_{\pi}[S(X, \mathbf{U}) q(X, \mathbf{U})].\end{aligned}$$

¹⁹Yes, I've abused notation, but it makes things cleaner... bear with me...

Gradient Factorisation

Recall Sutton's vanilla policy gradient:¹⁹

$$\begin{aligned}\nabla_{\theta} J(\pi) &= \mathbb{E}_{\pi}[\nabla_{\theta} \log \pi(\mathbf{U} | X) q(X, \mathbf{U})], \\ &\doteq \mathbb{E}_{\pi}[S(X, \mathbf{U}) q(X, \mathbf{U})].\end{aligned}$$

In this case we have the following properties:

- ▶ The target $q(X, \mathbf{U})$ is a scalar.
- ▶ The scores $S(X, \mathbf{U})$ are an $|\theta| \times 1$ real vector.

¹⁹Yes, I've abused notation, but it makes things cleaner... bear with me...

Gradient Factorisation

The **factored policy gradient** naturally extends this:

$$\nabla_{\theta} J(\pi) = \mathbb{E}_{\pi}[S(X, U) K \lambda \circ q(X, U)].$$

Gradient Factorisation

The **factored policy gradient** naturally extends this:

$$\nabla_{\theta} J(\pi) = \mathbb{E}_{\pi}[S(X, U) K \lambda \circ q(X, U)].$$

In this case we have the following properties:

- ▶ The target $\lambda \circ q(X, U)$ is an $m \times 1$ real vector.
- ▶ The scores $S(X, U)$ are an $|\theta| \times n$ real matrix.
- ▶ The influence K is an $n \times m$ real matrix.

Gradient Factorisation

The **factored policy gradient** naturally extends this:

$$\nabla_{\theta} J(\pi) = \mathbb{E}_{\pi}[S(X, \mathbf{U}) \mathbf{K} \circ \boldsymbol{\lambda} q(X, \mathbf{U})].$$

These are *equivalent in expectation due to the baseline trick*, with:

$$b_i(x, \mathbf{u}) \doteq [(\mathbf{1} - \mathbf{K}) \boldsymbol{\lambda} \circ q(x, \mathbf{u})]_i.$$

Market Making

Recall the mark-to-market reward function:

$$R_{n+1} = \underbrace{\delta_n^- \Delta \Omega_n^-}_{\text{Ask}} + \underbrace{\delta_n^+ \Delta \Omega_n^+}_{\text{Bid}} + \underbrace{\Omega_{n+1} \Delta Z_n (1 - \beta \Omega_{n+1})}_{?}.$$

Market Making

Recall the mark-to-market reward function:

$$R_{n+1} = \underbrace{\delta_n^- \Delta \Omega_n^-}_{\text{Ask}} + \underbrace{\delta_n^+ \Delta \Omega_n^+}_{\text{Bid}} + \underbrace{\Omega_{n+1} \Delta Z_n (1 - \beta \Omega_{n+1})}_{?}.$$

We can of course unroll ? to further separate variables:

$$R_n = R_n^- + R_n^+ + R_n^0,$$

where

$$R_{n+1}^\pm = \Delta Q_n^\pm \cdot (\delta_n^\pm \pm \Delta Z_n) + \beta \cdot \Delta Q_n^\pm \cdot (\Delta Z_n)^2 \cdot (2Q_n \pm \Delta Q_n^\pm),$$

$$R_{n+1}^0 = Q_n \cdot \Delta Z_n - \beta \cdot (\Delta Z_n)^2 \cdot (Q_n^2 - 2 \cdot \Delta Q_n^+ \cdot \Delta Q_n^-).$$

Market Making

We now have a clear split:

$$R_n = R_n^- + R_n^+ + R_n^0,$$

where

$$\begin{aligned}\text{corr}(R_{n+1}^+, R_{n+1}^- | X_n, U_n) &= 0, \\ \text{corr}(R_{n+1}^\pm, R_{n+1}^0 | X_n, U_n) &\neq 0.\end{aligned}$$

Market Making

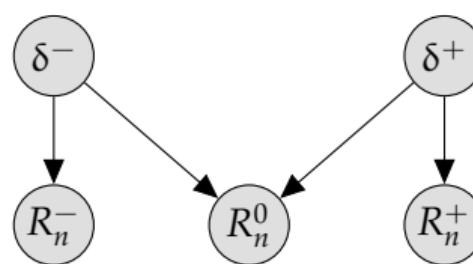
We now have a clear split:

$$R_n = R_n^- + R_n^+ + R_n^0,$$

where

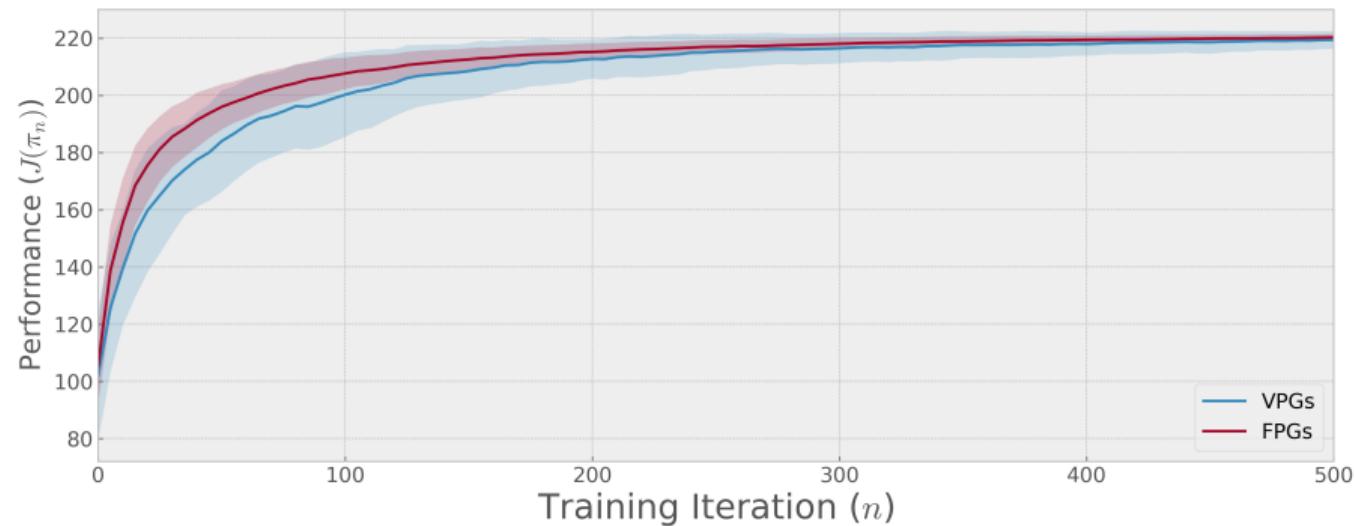
$$\begin{aligned} \text{corr}(R_{n+1}^+, R_{n+1}^- | X_n, U_n) &= 0, \\ \text{corr}(R_{n+1}^\pm, R_{n+1}^0 | X_n, U_n) &\neq 0. \end{aligned}$$

The resulting *influence network*:



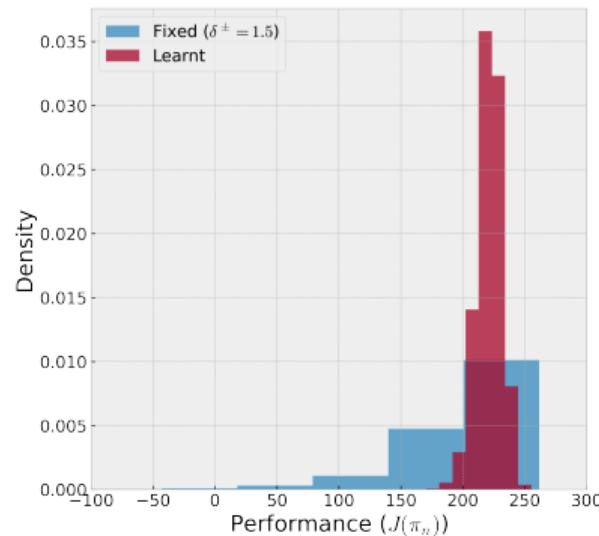
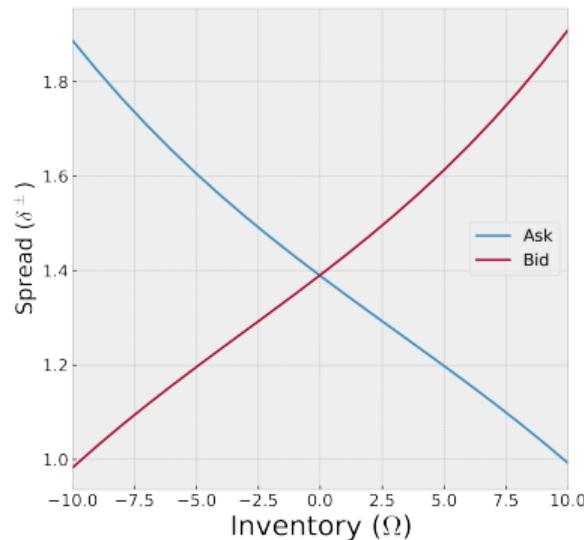
Market Making

Factoring gradients reduces variance!



Market Making

The (finite-time) policy is also more symmetric due to reduced variance!



Transition Acyclicity

How hard is it to solve MDPs in general?

Complexity

• Littman et al. (2013)



394

On the Complexity of Solving Markov Decision Problems

Michael L. Littman, Thomas L. Dean, Leslie Pack Kaelbling
 Department of Computer Science
 Brown University
 115 Waterman Street
 Providence, RI 02912, USA
 Email: {mlittman,tlid,lpk}@cs.brown.edu

Abstract

Markov decision problems (MDPs) provide the foundations for a number of problems of interest to AI researchers studying automated planning and reinforcement learning. In this paper, we present new results regarding the complexity of solving MDPs and the running time of MDP solution algorithms. We argue that, although MDPs can be solved efficiently in theory, more study is needed to find practical algorithms for solving large problems quickly. To encourage future research, we sketch some alternative methods of analysis that rely on the structure of MDPs.

1 INTRODUCTION

A *Markov decision process* is a controlled stochastic process satisfying the Markov property with costs assigned to state transitions. A *Markov decision problem* is a Markov decision process together with a performance criterion. A solution to a Markov decision problem is a policy, mapping states to actions, that (perhaps stochastically) determines state transitions to minimize the cost associated with the performance criterion. Markov decision problems (MDPs) provide the theoretical foundations for decision-theoretic planning, reinforcement learning, and other sequential decision-making tasks of interest to researchers and practitioners in computer science, operations research, and management science (Puterman, 1994). MDPs emphasize dynamic models based on well-understood stochastic processes and performance criteria based on established theory in operations research, economics, combinatorial optimization, and the social sciences (Puterman, 1994).

It would seem that MDPs exhibit special structure that might be exploited to expedite their solution. That is, interestingly, for many often used initial state is known with certainty (the current price for a stock or commodity) and as a result the set of likely reachable states (future prices) and viable invest-

ment strategies in the near-term future is considerably restricted. In general, notions of time, action, and reachability in state space are inherent characteristics of MDPs that might be exploited to produce efficient algorithms for solving them. It is important that we understand the computational issues involved in those sources of structure to get some idea of the prospects for efficient sequential and parallel algorithms for computing both exact and approximate solutions.

This paper summarizes some of what is known (and unknown but worth knowing) about the computational complexity of solving MDPs. In theory, any MDP can be represented as a linear program (LP) and solved in polynomial time. However, the order of the polynomials is large enough that the theoretically efficient algorithms are not efficient in practice. Of the algorithms known to solve MDPs, none is known to run in worst-case polynomial time. However, algorithms and analyses to date have made little use of MDP-specific structure, and results in related areas of Monte Carlo estimation and Markov chain theory suggest promising avenues for future research. We begin by describing the basic class of problems.

2 MARKOV DECISION PROBLEMS

For our purposes, a *Markov decision process* is a four-tuple $(\Omega_S, \Omega_A, p, c)$, where Ω_S is the state space, Ω_A is the action space, p is the state-transition probability-distribution function, and c is the instantaneous-cost function.

The state-transition function is defined as follows: for all $i, j \in \Omega_S, k \in \Omega_A$,

$$p_{ij} = \Pr(S_t = j | S_{t-1} = i, A_t = k)$$

where $S_t(A_t)$ is a random variable denoting the state (action) at time t . The cost c_t^k is defined to be the cost of taking action k at time t .

Let $N = |\Omega_S|$ and $M = |\Omega_A|$. For some of the computations, it will be necessary to assume that p and c are encoded using $N \times N \times M$ tables of rational numbers. We let B be the maximum number of bits required to represent any component of p or c . In

Complexity

- Littman et al. (2013)
- Fearnley and Savani (2014)

[arXiv:1404.0605v2 \[cs.DS\]](https://arxiv.org/abs/1404.0605v2) 17 Apr 2014

The Complexity of the Simplex Method

John Fearnley and Rahul Savani
University of Liverpool

Abstract. The simplex method is a well-studied and widely-used pivoting method for solving linear programs. When Dantzig originally formulated the simplex method, he gave a natural pivot rule that converts it into a polynomial-time algorithm with bounded cost. In their famous paper, Klee and Minty showed that this pivot rule takes exponential time in the worst case. We prove two main results on the simplex method. Firstly, we show that it is PSPACE-complete to find the solution that is computed by the simplex method using Dantzig's pivot rule. Secondly, we prove that deciding whether Dantzig's pivot rule finds a specific vertex in the simplex is NP-hard. This provides a new connection between Markov decision processes (MDPs) and linear programming, and an equivalence between Dantzig's pivot rule and a natural variant of policy iteration for average-reward MDPs. We contrast MDPs and show PSPACE-completeness results for single-switch policy iteration, which in turn imply our main results for the simplex method.

1 Introduction

Linear programming is a fundamental technique in computer science, and the simplex method is a widely used technique for solving linear programs. The simplex method requires a *pivot rule* that determines which variable is pivoted into the basis in each step. Dantzig's original formulation of the simplex method used a particularly natural pivot rule: in each step, the non-basic variable with the most negative reduced cost is chosen to enter the basis [5]. We will call this *Dantzig's pivot rule*. Klee and Minty have shown that Dantzig's pivot rule takes exponential time in the worst case [17].

The simplex method is a member of a much wider class of *local search* algorithms. The complexity class PLS, which was introduced by Johnson, Papadimitriou, and Yannakakis, captures problems where a locally optimal solution can be found by a local search algorithm [16]. PLS has matured into a robust class, and there is now a wide range of problems that are known to be PLS-complete [15, 19]. It is widely believed that PLS-complete problems do not admit polynomial-time algorithms.

To show that a problem lies in PLS, we must provide three polynomial-time functions: a function A that produces a candidate solution, a function B that assigns an value to each candidate solution, and a function C that for each candidate solution either produces a neighbouring candidate solution with higher value, or reports that no such candidate solution exists. Thus, each PLS problem comes equipped with a *natural algorithm* that executes A to find an initial candidate solution, and then iterates C until a local optimum is found. For some problems in PLS, it is known that it is PSPACE-complete to find the solution that is computed by the natural algorithm [22, 23]. So far, this is only known to hold for problems that are tight PLS-complete (see, e.g., [28]), which is a stronger form of PLS-completeness, or for problems that are suspected to be tight PLS-complete, such as the local max-cut problem on graphs of degree four [26].

Obviously, since linear programming is in P, it cannot be PLS-complete unless PLS=P. Despite this fact, in the first main theorem of this paper, we show that it is PSPACE-complete to compute the solution found by the simplex method equipped with Dantzig's pivot rule. Given a linear program \mathcal{L} , an initial basic feasible solution b , and a variable v , the problem $\text{DANTZIGLevSol}(\mathcal{L}, b, v)$ asks the following question: if Dantzig's pivot rule is started at basis b , and finds an optimal solution

Complexity

- Littman et al. (2013)
- Fearnley and Savani (2014)
- Balaji et al. (2018)

On the Complexity of Value Iteration

Nikhil Balaji
University of Oxford
nikhil.balaji@cs.ox.ac.uk

Stefan Kiefer
University of Oxford
stefkie@cs.ox.ac.uk

Petr Novotný
Masaryk University
petr.novotny@fjfi.zcu.cz
Guillermo A. Pérez
University of Antwerp
guillermo.alberto.perez@uaawpen.be

Mahsa Shirmohammadi
CNRS & IRIF
mahsa.shirmohammadi@irif.fr

Abstract

Value iteration is a fundamental algorithm for solving Markov Decision Processes (MDPs). It computes the maximal n -step payoff by iterating n times a recurrence equation which is naturally associated to the MDP. At the same time, value iteration provides a policy for the MDP that is optimal on a given set of states. In this paper, we study the computational complexity of value iteration. We show that, given a learned π in binary and an MDP, computing an optimal policy is EXPTIME-complete, thus resolving an open problem that goes back to the seminal 1967 paper on the complexity of MDPs by Papadimitriou and Tsitsiklis. To obtain this main result, we develop several stepping stones that yield results of an independent interest. For instance, we show that it is EXPTIME-complete to compute the n -fold iteration (with n is binary) of a function given by a straight-line program over the integers with max and $+$ as operators. We also provide new complexity results for the bounded halting problem in linear-update counter machines.

2012 ACM Subject Classification Theory of computation → Probabilistic computation; Theory of computation → Logic and verification; Theory of computation → Markov decision processes

Keywords and phrases Markov decision processes, Value iteration, Formal verification

Funding. Stefan Kiefer is supported by a Royal Society University Research Fellowship. Petr Novotný is supported by the Czech Science Foundation, grant no. G11-19-01349Y "Verification and Analysis of Probabilistic Programs." Mahsa Shirmohammadi is supported by the "AAPPS" PEPS JCJC grant.

Acknowledgements We thank James Worrell for helpful comments on early version of this work.

1 Introduction

Markov decision processes (MDP) are a fundamental formalism of decision making under probabilistic uncertainty [29, 30]. As such, they play a prominent role in numerous domains, including artificial intelligence and machine learning [14, 33], control theory [30, 3], operations research and finance [11, 33], as well as formal verification [12, 2], to name a few. Informally, an MDP represents a system which is, at every time step, in one of the states from a finite set S . The system evolves in steps: in each step, we can perform an action (or decision) from a finite set A . When using an action $a \in A$ in state $s \in S$, we collect an immediate reward

arXiv:1807.04920v3 [cs.FL] 27 Apr 2019

Complexity: Cycles

Consider the following undiscounted MRP:

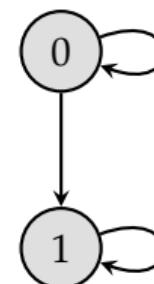
$$P \doteq \begin{bmatrix} \nu & 1 - \nu \\ 0 & 1 \end{bmatrix} \quad \text{and} \quad R \doteq \begin{bmatrix} \rho & 0 \\ 0 & 0 \end{bmatrix}.$$

We can of course solve this analytically...

$$v(0) = \underbrace{\nu [\rho + v(0)]}_{\text{Cycle}} + \underbrace{(1 - \nu) [0 + v(1)]}_{\text{Termination}},$$

such that

$$v(0) = \frac{\nu\rho}{1 - \nu}.$$



Complexity: Cycles

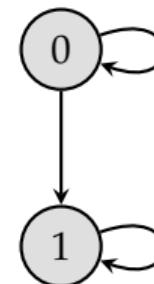
But how does value iteration perform?

Starting with $v_0(0) = 0$, we get

$$v_1(0) = \gamma\rho, \quad v_2(0) = \gamma\rho + \gamma^2\rho,$$

and eventually

$$v_k(0) = \sum_{i=1}^k \gamma^i \rho.$$



Complexity: Cycles

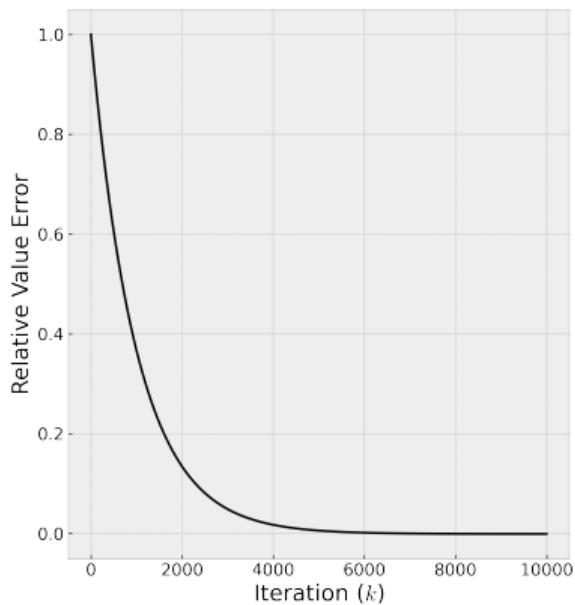
But how does value iteration perform?

Of course, we end up in the same place...

$$\lim_{k \rightarrow \infty} v_k(0) = \frac{\gamma \rho}{1 - \gamma}.$$

But it will take a long time to get there!

► Cycles are bad!



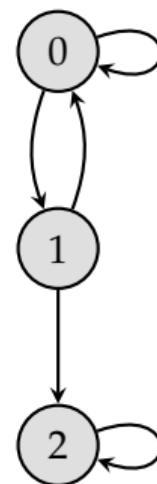
Complexity: Bootstrapping

But there's a bigger issue!

- ▶ *VI relies heavily on bootstrapping.*

When estimating $v(\cdot)$ for multiple states:

- ▶ Cycles get longer.
- ▶ Feedback is slow to propagate.
- ▶ And intermediate estimates are erroneous.



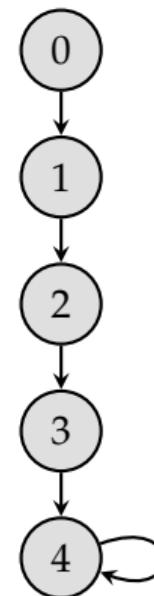
Acyclicity

What if there weren't any cycles?

- ▶ Could we do better?

The key is hidden in the VI algorithm:

- ▶ *We are free to choose the order of updates!*



Acyclicity

Theorem (Optimal Backup Order²⁰)

If an MDP is acyclic, then there exists an optimal backup order. By applying the optimal order, the optimal value function can be found with each state needing only one backup.

²⁰Dimitri Bertsekas. *Dynamic Programming and Optimal Control*. Vol. 2. Athena Scientific, 2012.

Acyclicity

Theorem (Optimal Backup Order²¹)

If an MDP is acyclic, then there exists an optimal backup order. By applying the optimal order, the optimal value function can be found with each state needing only one backup.

²¹Dimitri Bertsekas. *Dynamic Programming and Optimal Control*. Vol. 2. Athena Scientific, 2012.

Acyclicity

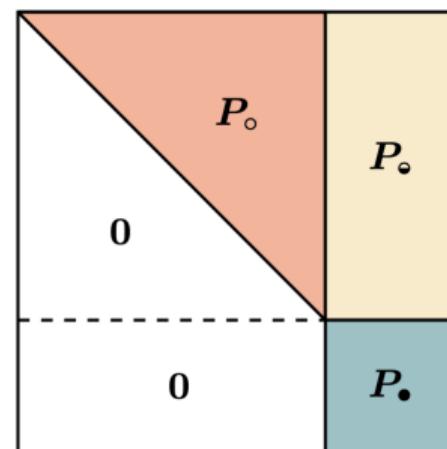
Why is this true?

- ▶ Acyclicity implies P is upper-triangular!

And we can take advantage of this:

- ▶ Perform VI updates according to the topological order.
- ▶ “Start from the end, and work backwards.”

Recall the backwards induction algorithm!



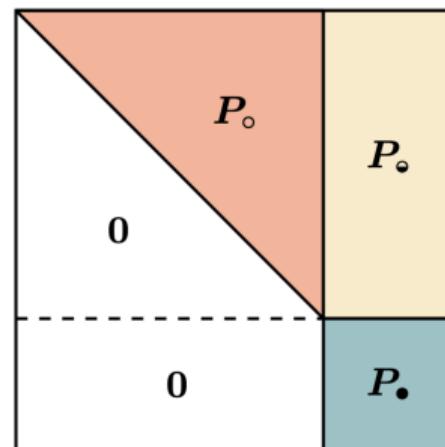
Acyclicity

It turns out, many problems are like this!

1. Routing problems.
2. Finite resource problems.
3. *Optimal liquidation.*

In other cases, we can impose it.

- ▶ This leads to sub-optimality.
- ▶ But drastically reduces training time.



Acyclicity: Optimal Liquidation

Let's consider a concrete example.

- ▶ Discrete optimal liquidation problem.
- ▶ *Where the agent has to trade every step.*

Acyclicity: Optimal Liquidation

Let's consider a concrete example.

- ▶ Discrete optimal liquidation problem.
- ▶ *Where the agent has to trade every step.*

Inventory Process

- ▶ Agent chooses $0 < U_n \leq \Omega_n$.
- ▶ Inventory transitions:

$$\Omega_n \mapsto \Omega_n - U_n.$$

Acyclicity: Optimal Liquidation

Let's consider a concrete example.

- ▶ Discrete optimal liquidation problem.
- ▶ *Where the agent has to trade every step.*

Inventory Process

- ▶ Agent chooses $0 < U_n \leq \Omega_n$.
- ▶ Inventory transitions:

$$\Omega_n \mapsto \Omega_n - U_n.$$

Price Process

- ▶ Trinomial random walk: $p \in \Delta^2$.
- ▶ With imposed reflection.

Acyclicity: Optimal Liquidation

Let's consider a concrete example.

- Discrete optimal liquidation problem.
- *Where the agent trades every step.*

Reward Process

$$R_n \doteq \underbrace{w_0 \cdot U_n \cdot (Z_n - Z_0)}_{\text{Excess PnL}} - \underbrace{w_1 \cdot U_n^2 - w_2 \cdot \Omega_n^2}_{\text{Penalty}}$$

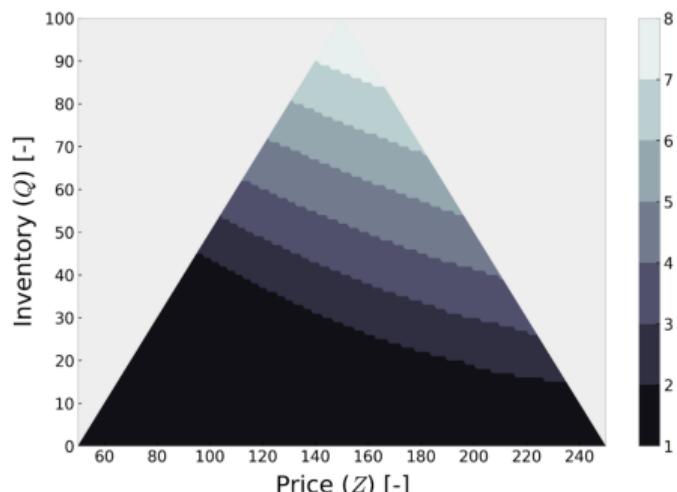
Acyclicity: Optimal Liquidation

Let us set the dynamics to

- ▶ Start with $\Omega_0 \doteq 100$ and $Z_0 \doteq 150$.
- ▶ Transitions with $p \doteq [0.4, 0.2, 0.4]$.
- ▶ And reward weights $w \doteq [1, 0.2, 0.002]$.

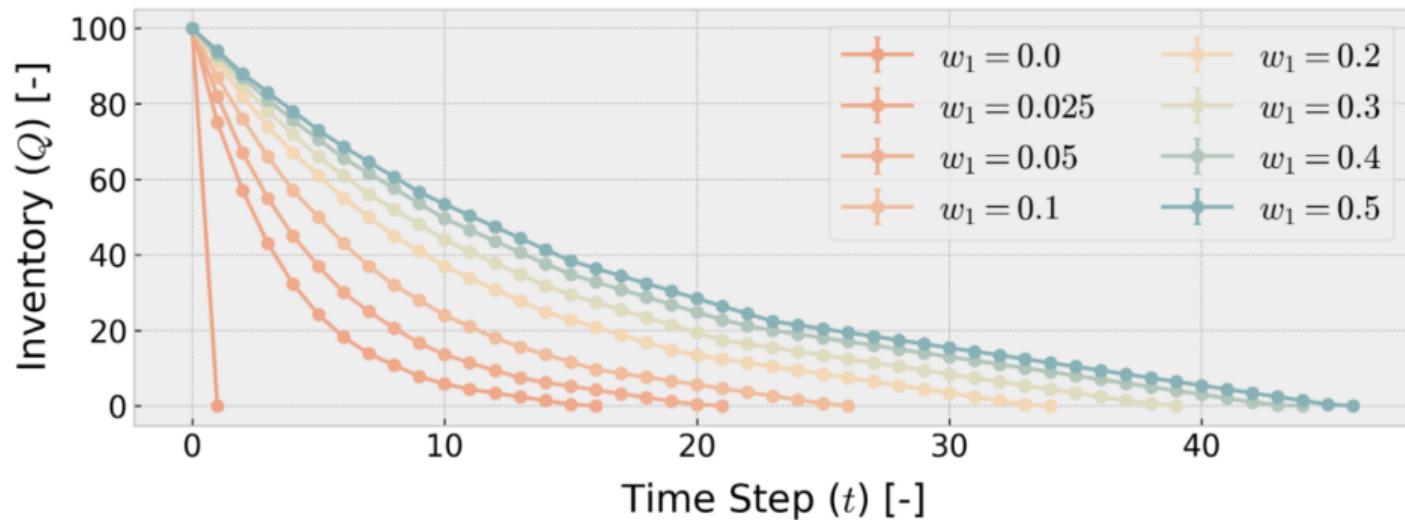
Optimal policy is pretty intuitive!

- ▶ Sell fast when prices are high.
- ▶ Decrease rate as inventory shrinks.



Reductivity

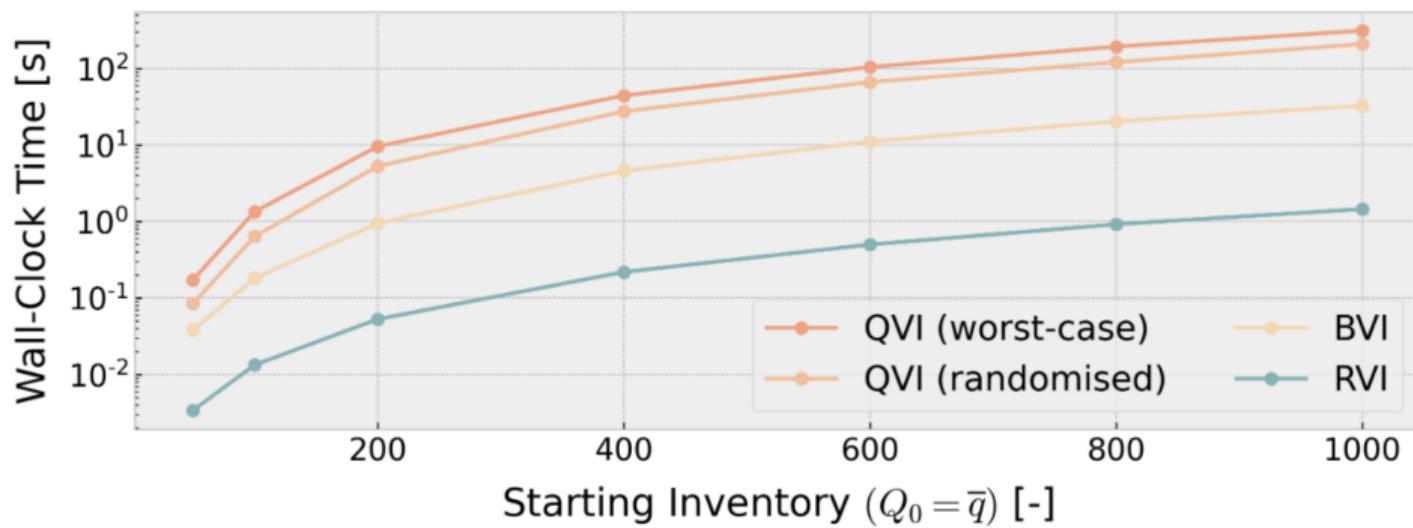
Increasing the penalty term w_1 slows the trading rate.



Reducitivity

Accounting for acyclicity (RVI) gives OoM reductions in time!

- ▶ Time complexity of RVI forms a lower-bound.²²



²²Leveraging structure for fun and profit...and saving time!

Reductivity

So what's the theory behind it?

Reductivity

Recall: a set $B \in \mathfrak{B}(\mathsf{X})$ is reachable from $x \in \mathsf{X}$ if $P_n(x, B)$ for some $n > 0$.

We thus define the *n-step reachable state set* from $x \in \mathsf{X}$ as

$$\mathsf{F}_0(x) \doteq \{x\},$$

$$\mathsf{F}_n(x) \doteq \bigcap \{B \in \mathfrak{B}(\mathsf{X}) : P_n(x, B) > 0\}.$$

The *reachable state set* is the union thereof

$$\mathsf{F}(x) \doteq \bigcup_{n=0}^{\infty} \mathsf{F}_n(x) \quad \implies \quad \mathsf{F}(B) \doteq \bigcup_{x \in B} \mathsf{F}(x).$$

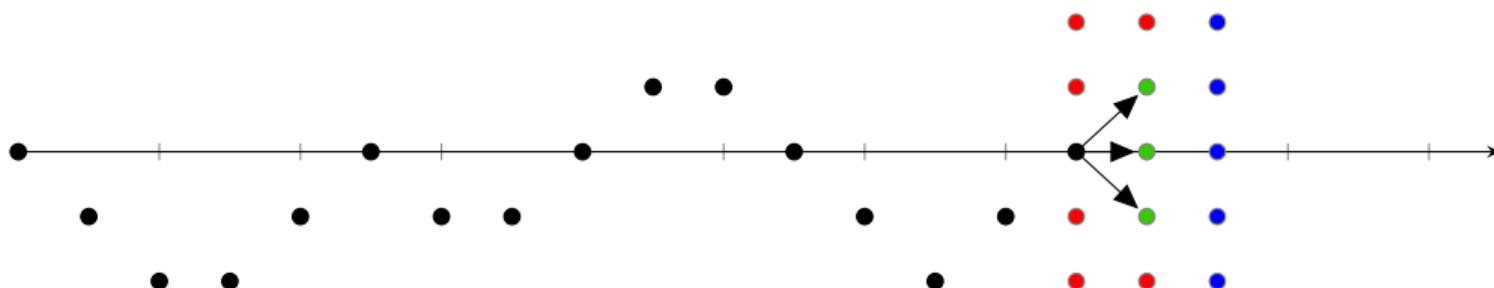
Reductivity

In the random walk, where $P(x, \{x + 1\}) = P(x, \{x\}) = P(x, \{x - 1\})$, we have

$$\mathsf{F}_1(x) = \{x - 1, x, x + 1\},$$

$$\mathsf{F}_2(x) = \{x - 2, x - 1, x, x + 1, x + 2\},$$

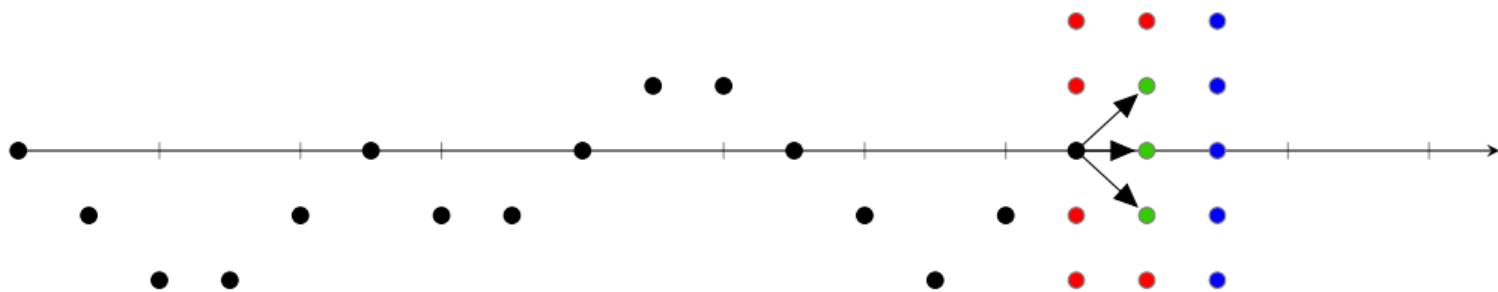
⋮



Reductivity

In the random walk, where $P(x, \{x + 1\}) = P(x, \{x\}) = P(x, \{x - 1\})$, we have

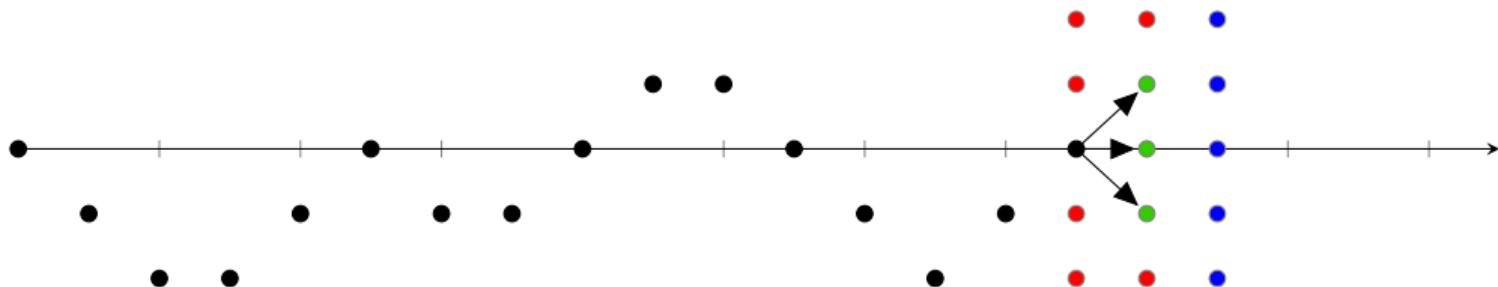
$$\mathsf{F}_n(x) = \{x - n, \dots, x, \dots, x + n\}.$$



Reductivity

In the random walk, where $P(x, \{x + 1\}) = P(x, \{x\}) = P(x, \{x - 1\})$, we have

$$\mathsf{F}(x) = \bigcup_{n=0}^{\infty} \{x - n, \dots, x + n\} = \mathbb{Z}.$$



Reductivity

Is there a succinct way to describe reachable sets?

Reductivity

Is there a succinct way to describe reachable sets? Yes!

For a measure $\mu : \mathcal{B}(X) \rightarrow [0, \infty)$, the *reachability potential* is

$$\phi_\mu(x) \doteq \mu \circ F(x),$$

with “discrete derivative” given by

$$\psi_\mu(x, x') \doteq \phi_\mu(x') - \phi_\mu(x) = -\mu(F(x) \setminus F(x')) \leq 0.$$

Reductivity

If we now define the set of non-absorbing cyclic states,

$$\mathsf{L}(A) \doteq \{x \in A : 0 < \mathbb{P}(x, \{x\}) < 1\},$$

then we can now define an RMC.

Reductivity

If we now define the set of non-absorbing cyclic states,

$$\mathsf{L}(A) \doteq \{x \in A : 0 < \mathbb{P}(x, \{x\}) < 1\},$$

then we can now define an RMC.

Definition (RMC)

Let μ be a measure. Then, an MC is said to be μ -reductive if $X_\bullet \neq \{\}$ and

$$\begin{cases} \psi_\mu(x, x') = 0 & \text{for } x' \in X_\bullet \cup \mathsf{L}(\{x\}), \\ \psi_\mu(x, x') < 0 & \text{otherwise,} \end{cases}$$

holds for all states $x \in X$ and their successors $x' \in F_1(x)$.

Reductivity

This is a sufficient condition for acyclicity!

- ▶ But also extends the concept to continuous state-spaces!

Definition (RMC)

Let μ be a measure. Then, an MC is said to be μ -reductive if $X_\bullet \neq \{\}$ and

$$\begin{cases} \psi_\mu(x, x') = 0 & \text{for } x' \in X_\bullet \cup L(\{x\}), \\ \psi_\mu(x, x') < 0 & \text{otherwise,} \end{cases}$$

holds for all states $x \in X$ and their successors $x' \in F_1(x)$.

Reductivity

We can also apply to MDPs (see below):²³

- *The corresponding algorithm, RVI, solves RMDPs in polynomial time!*

Definition (RMDP)

An MDP is said to be μ -reductive if every policy induces a μ -RMC.

²³Thomas Spooner et al. “Reductive MDPs: A Perspective Beyond Temporal Horizons”. In: arXiv preprint arXiv:2205.07338 (2022).

That's all folks!