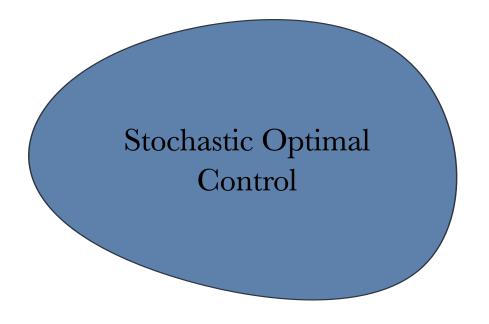
Connecting SOC with RL – Importance sampling

Alonso Cisneros

Zuse Institute Berlin

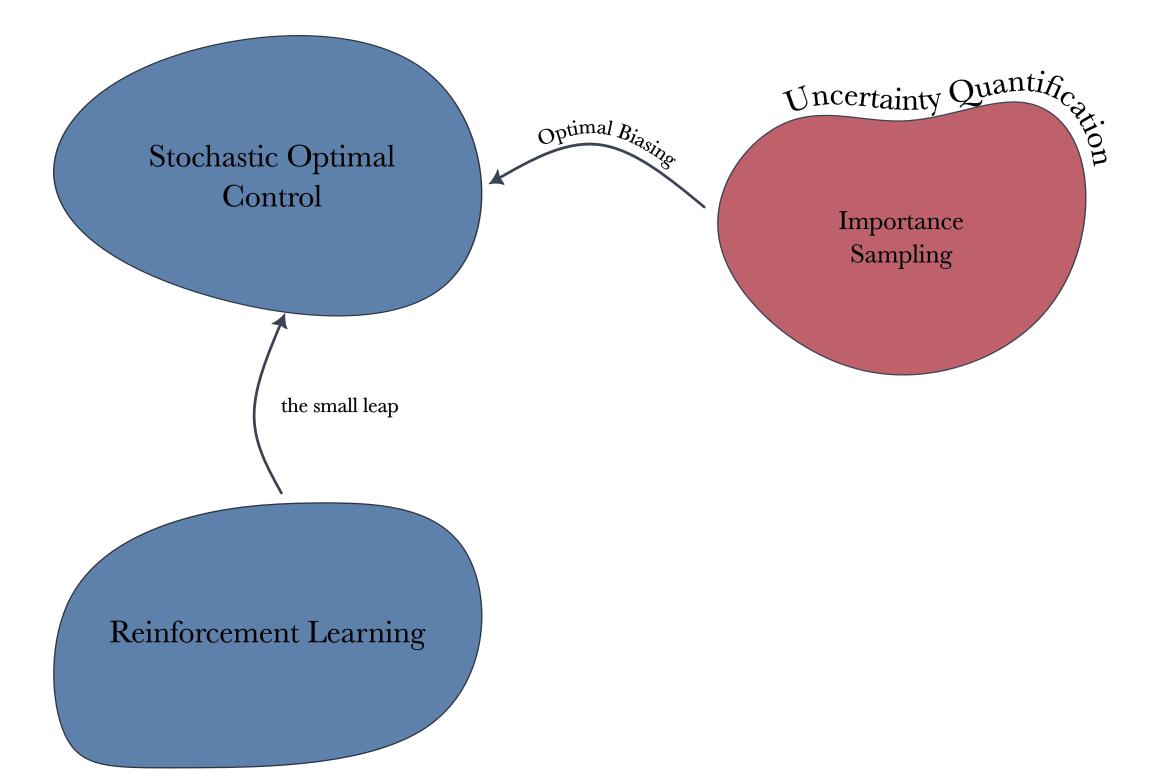


Uncertainty Quantificants

Importance
Sampling

Reinforcement Learning

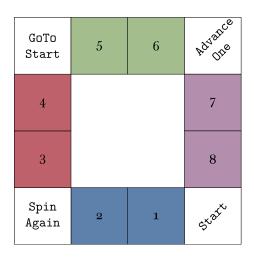
By the end of the talk these will all be connected. What I want to do is build the connection to yet another blob: ABMs



Outline

- 1. Crash course on RL
- 2. What is importance sampling
 - The connection to optimization
 - Optimal biasing
- 4. Optimal biasing as an RL problem
- 5. The things I'd like to connect

Crash course on Reinforcement Learning



A miniopoly board

- I'm training a robot to become the best Miniopoly player
- The rules:
 - Players play in turns
 - They move a number of squares determined by a 4-sided dice roll
 - \blacksquare After completing a lap, it gets a reward of x dollars
 - The trap squares do what it says on the square
 - They can buy property and hotels in the squares.
 - If they land of a square someone owns, they pay
 - o If someone lands on their square, they charge rent
 - The game ends when someone runs out of money

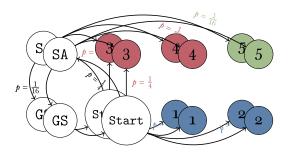
- ullet The game has a state at turn t denoted s_t
- At a turn t players roll the dice
- ullet The change in money after buying/paying rent/charging rent is recorded as a reward r_t

! Important

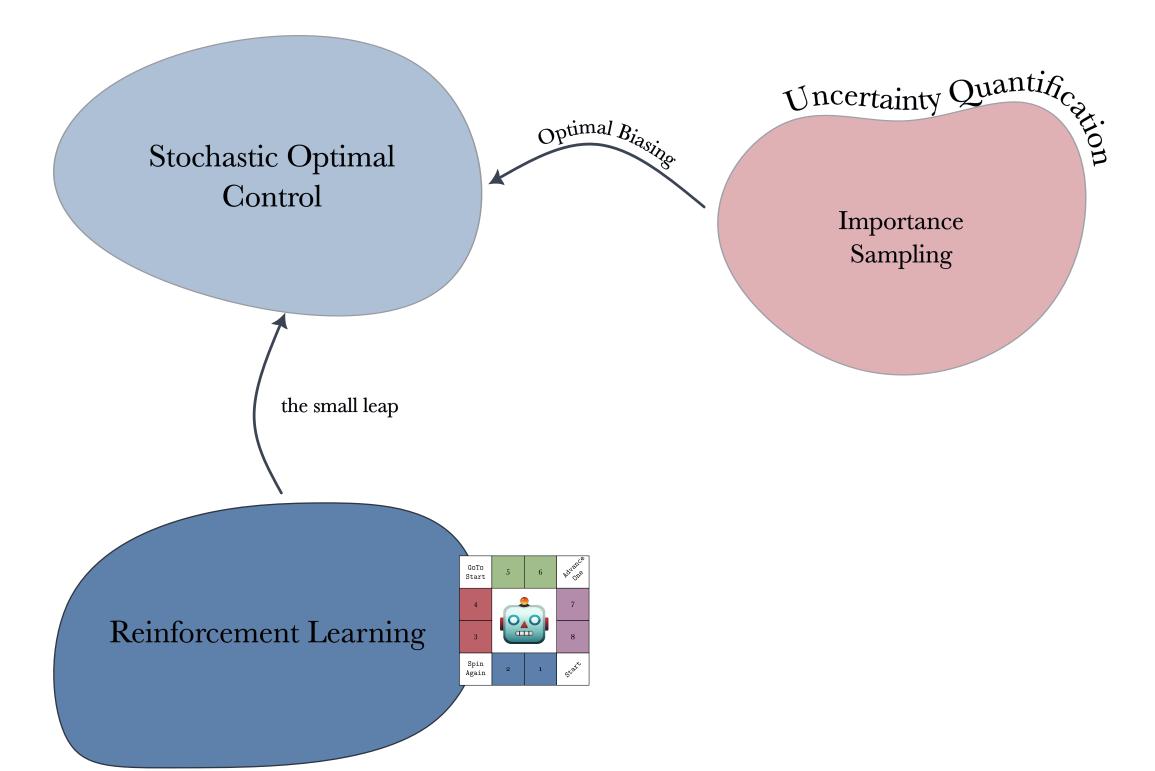
We train our robot to maximize the rewards as it takes actions exploring the space of states

- The state of the game an any given time is information like, who owns what squares, how much money they have, in what positions each player is, and so on.
- Once a player lands in another square, they can choose to buy it if available. If it's not, we carry out the accounting of how much rent is, and let the player know how much it won/lost and to what square this is connected.

Dynamics of the space



In this example we have full access to the dynamics of the problem. This is not always the case



Speaker notes • We just reviewed very quickly the first bubble

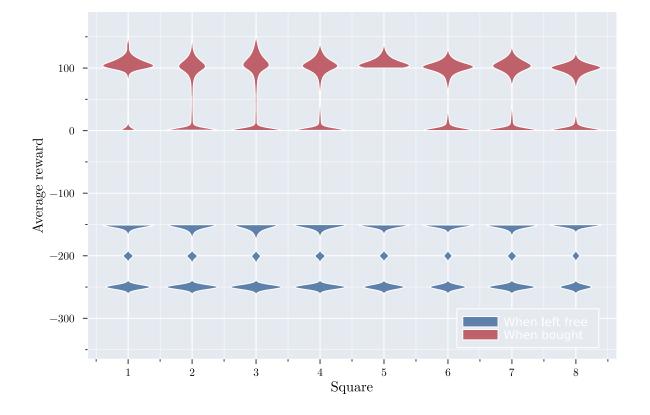
What if we don't know how square transitions work?

We calculated transition probability with the knowledge of the dice

We'll now move on to MCMC. These concepts are not necessarily related, but I will take advantage of this example so that we can have a reasonable inutition from the start.

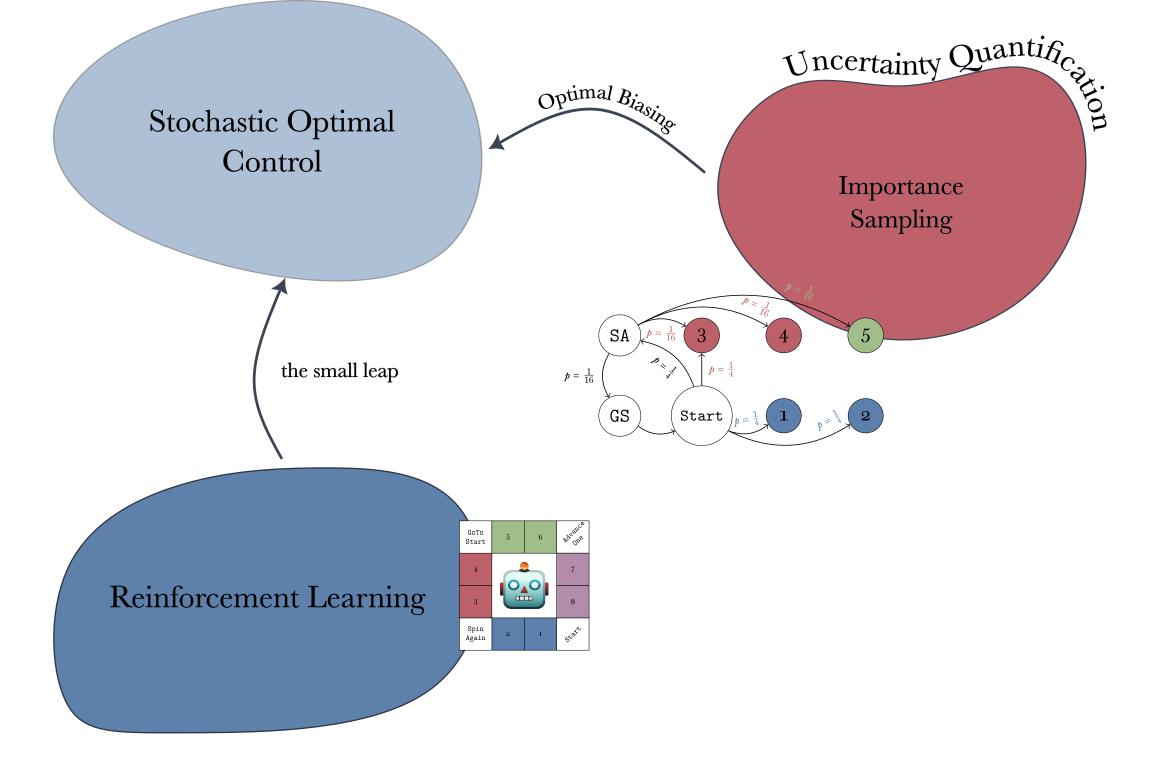
Markov Chain Monte Carlo

- We let the robot roam around and buy squares as it pleases
 - For any square, it can either buy it or not
- We register what it gained or lost by buying or not buying a square by the end of the game.



• Graph:

- This is a violin plot. It shows the estimated probability densities of observing...
- The x axis shows the different squares
- The y axis shows the estimated rewards. i.e. money
- The red distributions correspond to the expected reward when buying a square, while the blue the expected loss when not buying them
- i.e. When we buy squares we expect to profit from them, but clearly not all squares are as profitable, look at the different shapes of the distributions. On the other hand, it looks like losing any given square leads to the same expected loss



Spe What we've covered so far

Now we have introduced MCMC as a way for our robot to explore it's environment and estimate which moves are beneficial to his goal. This is what the diagram is representing

Moving on...

Importance Sampling

- We wanted to compute the expected reward of the robot after the entire game
- MCMC often fails in metastable systems
- Importance sampling aims to remedy this

! Important

The general idea of importance sampling is to draw random variables from another probability measure and subsequently weight them back in order to still have an unbiased estimator of the desired quantity of interest

- We **estimated** this quantity by observing and measuring an empirical average. But our approximation for extremely unlikely states will always be bad by virtue of how little samples we have.
- Metastability makes MCMC extremely hard to apply. The variance of our estimations is always going to be enormous under these conditions
- We can aim to make sampling faster by reducing the inherent variance
- After Callout In the case of stochastic processes this change of measure corresponds to adding a control to the original process

More formally...

- Where:
 - $lacksquare X_s$ is the position of our particle at time s
 - *V* is a "potential"
 - We assume there exists a unique strong solution that is ergodic
- Note that τ is a.s. finite
- ullet Where ${\mathcal W}$ serving as a measure of "work" over a trajectory

Our main goal is to compute

$$\Psi(X)\coloneqq \mathbb{E}^x[I(X)]\coloneqq \mathbb{E}[I(X)\mid X_0=x]$$

But...

- MCMC has terrible properties because of metastability
- ullet Closed forms of $\Psi(X)$ maybe don't exist

○ Tip

- We can "push" the particle adding force, as long as we account for it and correct for that bias
- That "push" is achieved by adding a control u.

The new, controlled dynamics are now described as

$$\mathrm{d}X^u_s = (-\nabla V(X^u_s) + \sigma(X^u_s)\,u(X^u_s))\mathrm{d}s + \sigma(X^u_s)\mathrm{d}W_s,$$

Via Girsanov, we can relate our QoI to the original as such:

$$\mathbb{E}^x\left[I(X)
ight]=\mathbb{E}^x\left[I(X^u)M^u
ight],$$

where the exponential martingale

$$M^u \coloneqq \exp\left(-\int_0^{ au^u} u(X^u_s)\cdot \mathrm{d}W_s - rac{1}{2}\int_0^{ au^u} |u(X^u_s)|^2 \mathrm{d}s
ight)$$

corrects for the bias the pushing introduces.

! Important

The previous relationship always holds. But the variance of the estimator depends heavily on the choice of u.

Clearly, we aim to achieve the smallest possible variance through on optimal control u^{\ast}

$$\operatorname{Var}\left(I(X^{u^*})M^{u^*}
ight) = \inf_{u \in \mathcal{U}} \left\{\operatorname{Var}(I(X^u)M^u)
ight\}$$

- Where:
 - $lacksquare X^u_s$ is the position of our particle at time s under control u
 - lacktriangle The potential u is an Itô integrable function satisfying a linear growth condition
- Note that τ is a.s. finite
- ullet Where ${\mathcal W}$ serving as a measure of "work" over a trajectory

Connection to optimization

It turns out ¹ that the problem of minimizing variance corresponds to a problem in optimal control

The cost functional J to find the variance minimizing control is

$$J(u;x)\coloneqq \mathbb{E}^x\left[\mathcal{W}(X^u)+rac{1}{2}\int_0^{ au^u}|u(X^u_s)|^2\mathrm{d}s
ight],$$

With this formulation,

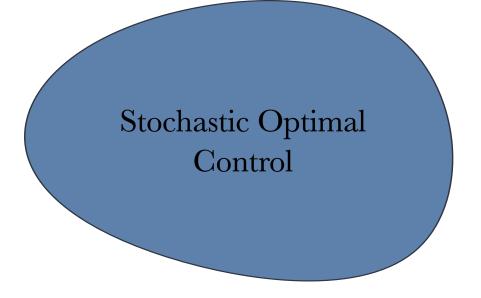
$$\Phi(x) = \inf_{u \in \mathcal{U}} J(u;x).$$

! Important

The optimal bias achieves zero variance:

$$\operatorname{Var}\left(I(X^{u^*})M^{u^*}
ight)=0.$$

Optimal biasing through RL



Uncertainty Quantificants

Importance
Sampling

Reinforcement Learning

- Let's reconsider the SOC problem (excuse the change in notation)
- We discretize the dynamics equation

$$egin{aligned} s_{t+1} &= s_t + (-
abla V(s_t) + \sigma u(s_t))\Delta t + \sigma \sqrt{\Delta t}\,\eta_{t+1} \ s_0 &= x \end{aligned}$$

- Sorry for the slightly different notation
- Where
 - lacktriangle Our state is now represented by s
 - $\quad \blacksquare \ \ \text{We have the same potential } V$
 - lacktriangle The diffusion term is σ again
 - lacksquare Δt is the length of the time step
 - lacksquare The term $\sqrt{\Delta t}\eta_{t+1}$ is a Brownian increment, $\eta_t \sim N(0,1)$

The time-discretized objective function is given by

$$J(u;x) \coloneqq \mathbb{E}^x \left[g(s_{T_u}) + \sum_{t=0}^{T_{u-1}} f(s_t) \Delta t + rac{1}{2} \sum_{t=0}^{T_{u-1}} |u(s_t)|^2 \Delta t
ight]$$

ullet Our stopping time au is now denoted T_u

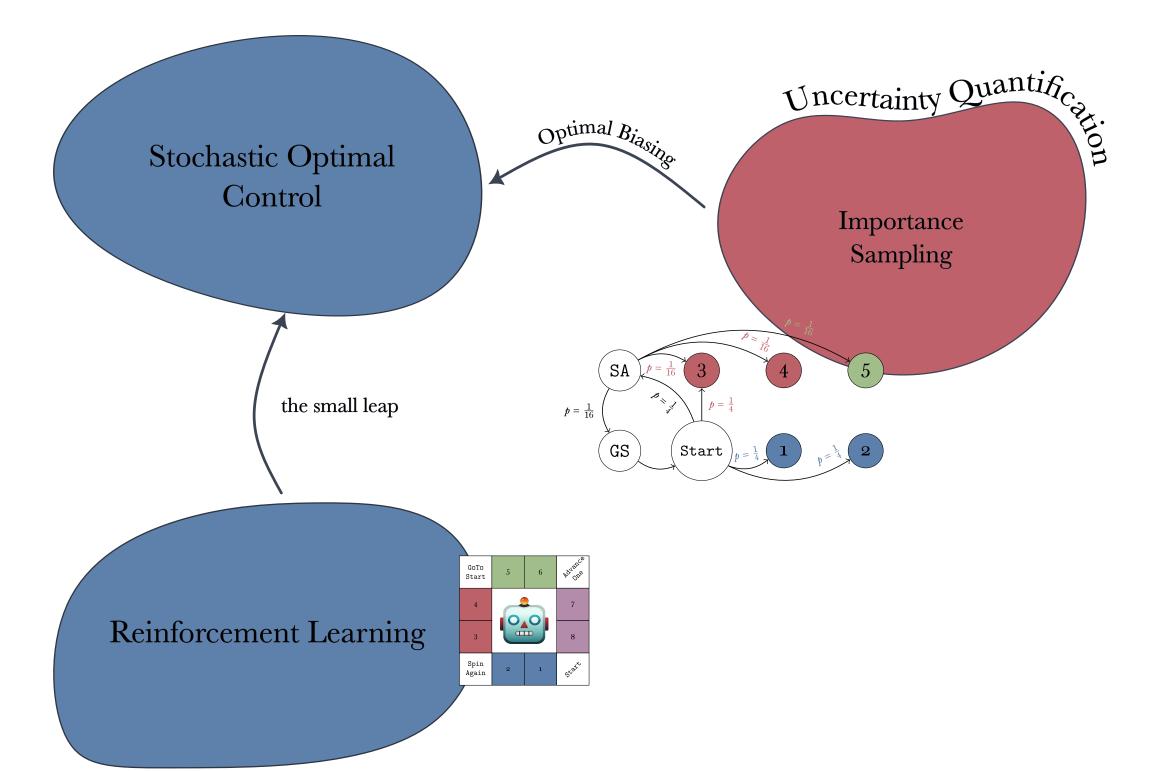
Some formalities

- ullet The state space ${\mathcal S}$ is all possible $s\in {\mathbb R}^d$
- ullet The action space ${\mathcal A}$ is the codomain of all possible controls ${\mathbb R}^d$
- The stopping time T_u for the controlled process is a.s. finite
- We'll approximate the control with Galerkin projections $u_{ heta}$
- We still need to derive probability transition and reward functions

The return we want to optimize depends on a rewards function

$$r_t = r(s_t, a_t) \coloneqq egin{cases} -f(s_t)\Delta t - rac{1}{2}|a_t|^2\Delta t & ext{if } s_t
otin \mathcal{T} \ -g(s_t) & ext{if } s_t
otin \mathcal{T} \end{cases}$$

- The reward function is defined such that the corresponding return along a trajectory equals the negative term inside the expectation of the time-discretized cost functional
- Notice that the reward signal is in general not sparse since the agent receives feedback at each time step but the choice of the running cost f and the final cost g can influence this statement.



What I want to do

- \bullet The connection (Quer and Borrell 2024) works because of the properties of J
- Two posibilities
 - 1. Take the SOC already published in (Helfmann et al. 2023) and pose the *right* cost functional (the "easy one")
 - 2. Go back to the MFE and then
 - Break the assumption of fully connected Agent-Agent network
 - Find the MFE without that assumption
 - Pose a SOC
 - Find the right cost functional

• Half of the magic is the fact that the functional to be optimized to solve the optimal biasing SOC problem leads to HJB. It's an open question whether I can make our SOC compatible

What I'm trying at the moment

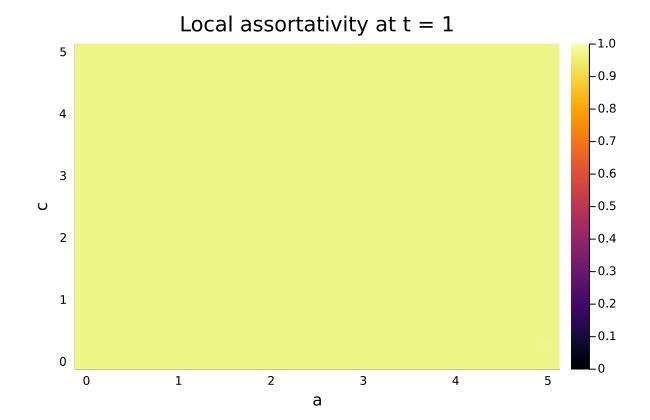
- Approach it from the MFE side
 - Looking for generalizations of McKean-Vlasov PDEs with time-dependent networks
 - Thinking what graphon I can get from the network
- Finding the right cost functional on the "easy" case

Bonus: Assortativity & network topology

Start

End

End (echo chamber)



References

- Helfmann, Luzie, Nataša Djurdjevac Conrad, Philipp Lorenz-Spreen, and Christof Schütte. 2023. "Modelling Opinion Dynamics Under the Impact of Influencer and Media Strategies." *Scientific Reports* 13 (1): 19375. https://doi.org/10.1038/s41598-023-46187-9.
- Quer, J., and Enric Ribera Borrell. 2024. "Connecting Stochastic Optimal Control and Reinforcement Learning." *Journal of Mathematical Physics* 65. https://doi.org/10.1063/5.0140665.
- Ribera Borrell, Enric, Jannes Quer, Lorenz Richter, and Christof Schütte. 2024. "Improving Control Based Importance Sampling Strategies for Metastable Diffusions via Adapted Metadynamics." *SIAM Journal on Scientific Computing* 46 (2): S298–323. https://doi.org/10.1137/22M1503464.