

Deep Reinforcement Learning in Macroeconomic Models

Matias Covarrubias

July 16, 2021

Abstract

Can artificial intelligent (AI) algorithms learn optimal intertemporal behavior without knowing any mathematical details of the economy beforehand? How should we specify economies and markets so as to make them easy to learn? In this paper, I apply Deep Reinforcement Learning (Deep RL), the current front-runner in the design of AI agents, to macroeconomic problems. In the context of an investment under uncertainty model, I show that Deep RL agents can learn the rational expectations solution and I suggest design choices that aid such learning. Then, by manipulating our baseline framework I highlight three topics in which this technology can be useful: high dimensional problems, incomplete information problems and equilibrium selection in models with multiple equilibrium.

1 Introduction

- Paragraph 1 (reframe) :Can artificial intelligent (AI) algorithms learn optimal intertemporal behavior without knowing any mathematical details of the economy beforehand? How should we specify economies and markets so as to make them easy to learn? In this paper, I apply Deep Reinforcement Learning (Deep RL), the current front-runner in the design of AI agents, to macroeconomic problems. In the context of an investment under uncertainty model, I show that Deep RL agents can learn the rational expectations solution and I suggest design choices that aid such learning. Then, by manipulating our baseline framework I highlight three topics in which this technology can be useful: high dimensional problems, incomplete information problems and equilibrium selection in models with multiple equilibrium.
- Paragraph 2: Explain model free RL.
 - RL is a class of algorithms that adapt dynamic programming techniques such as value function and policy function iteration to the problem of online learning, that is, to learn how to control a system by interacting and experimenting with it. Other historical names for this method are approximate dynamic programming, learning based control and automated optimization.
 - A key characteristic of RL algorithms is that they do not use any mathematical knowledge of the system they are controlling, other than the allowed actions and states. They feed actions to the system and observe the reward they get that period and how the state evolves. After millions of such transitions, if learning is successful, they are able to estimate the consequences of their actions accurately and thus make optimal decisions.
- Paragraph 3: Deep RL and the advent of this technology.
 - Deep RL adds neural nets as function approximator. Explain paper.
 - Why it has been so effective? Scalability and computation power.
 - Success cases: games, robotics, self-driving vehicles, industrial optimization.
- Paragraph 3: Traditional challenges of Deep RL and why economics may be particularly well suited for this technology.
 - Rewards engineering is usually the main bottleneck. In many problems, such as games or robotics, you only get positive rewards after successfully achieving some final objective, so the researcher needs to design signals to indicate to the agents whether they are making progress or not. In economics, rewards are well specified in every period.

- Also, RL algorithms usually require millions of transitions in order to learn, which may be computationally expensive. This problem is very serious in environments with high dimensional states, such as images, and with complex transitions logic, as in the case of games or robotics. In economics, sampling a transition is usually very fast, as we are only dealing with vectors whose transitions laws can be expressed with simple mathematical functions.
 - Finally, in many problems the algorithms get stuck in local optima. While this is also a potential problem in complex economic models, for many models we have nice regularity properties.
- Paragraph 5: But, economic problems present their own challenges.
 - Markets are a multi-agent problem. While multi agent Deep RL is a very active topic with some impressive success cases, it faces a crucial challenge: agents learn and explore while other agents learn and explore. This means that the environment is non-stationary and at the beginning it may be very chaotic. Whether order arises from such chaos is an empirical question with no theoretical guarantees.
 - In economics, rigorous insights often rely on exact solutions. That is, we are trying to solve very precise control problems, so usual assumptions such as discrete action spaces may not be satisfactory. Deep RL agents might get a very high percentage of the optimal discounted utility, but the solution may still lack some of the qualitative features of the exact solution. Example: growth model with fixed saving rate instead of decreasing saving rate.
 - Paragraph 4: Explain why this could be useful besides the study of learning in macroeconomics:
 - With this technology, we only need to specify the economy's action space, observation space, initial state and transition laws. This makes it simple to setup an economy on which Deep RL agents can learn.
 - The algorithm's structure is mostly the same regarding the problem at hand. The algorithm only needs to adapt the dimensions of its function approximators (e.g., neural nets) to the state space and action space of the environment, which can be easily automated. Thus, the same algorithm can learn to play chess or to calculate optimal saving behavior. As a corollary, the information available to each agent can be manipulated easily, without making any changes to the algorithm. Also, different versions of the economy can be constructed without need to recalculate the set of equations that describe the optimal solution

1.1 Preview of leading framework and exercises

- My criteria for choosing a baseline framework are the following:

- Simple economics, challenging solution: the problem should be easy to explain and with clear economics. Nevertheless, the control problem should be challenging.
 - Dynamics at the forefront: the key trade-offs of the model should be intertemporal, since we are interested in the application of deep RL to dynamic problems.
 - Meaningful scalability: the problem needs to have a clear path to scalability in which the dimensionality truly affects the economics of the problem. As such, it cannot have aggregation properties, so each variable in the state space should affect optimal behavior.
- Consistent with this criteria, I present a model of investment under uncertainty in which we focus on the purchase market for capital goods. There are a discrete number of buyers (household who produce, consume and invest) and sellers (each producing a different type of capital good). More precisely:
 - N^h households indexed by choose how much to consume and how much to invest in order to maximize the discounted intertemporal utility..
 - The final good is produced using N^c capital goods. All capital goods are bundled with a Cobb-Douglas aggregation and the production function have decreasing returns to scale over that bundle. The total factor productivity of each households follows a markov chain.
 - The agent allocates current production between consumption and investment in new capital goods. The cost in term of the final good of producing new capital goods is quadratic in investment.
 - We will consider both a market's formulation and a planner's formulation of the problem. In the market formulation, for each capital good there is a firm who pays the quadratic cost and sells the capital good. In the planner formulation, the planner decides how much of the final good to allocate in each capital good for each household, and investment is determined by how many units can be produced given the total allocation of resources by all households on that capital good.
 - Within this framework we perform the following exercise:
 - Exercise 1: The planner problem. For single agent problems, we can ask the Deep RL agents to solve the exact same problem we formulate in macroeconomic models. Thus, we can check how close does the Deep RL agent gets to the solution obtained by classical dynamic programming methods. If convergence is achieved, we confirm that the expectations implicit in the behavior of the Deep RL model are rational. Also, we can scale the problem to many state variables and evaluate the potential of this learning methodology to deal with high dimensional dynamic control problems.
 - Exercise 2: Many households.

- Exercise 3: Incomplete Information.
- Exercise 4: AI agents in both supply and demand.
- Paragraph 7: Explain learning results.
- Paragraph 8: Explain high dimensional, incomplete information, and multiple equilibrium exercises.
- Paragraph 9: Explain design lessons.

1.2 Related Literature

- paragraph 1: Learning in Macroeconomics.
- paragraph 2: Reinforcement Learning in economics.
- paragraph 3: Literature on machine learning learning for high-dimensional problems.

2 A Primer on Reinforcement Learning

- Define Markov Decision Problems.
- Example 1: Monopolist problem.
- Algorithm: Q-learning
 - We start from simplest single agent algorithm. The problem is:

$$\max_{a \in A} \sum_{t=s}^{\infty} \gamma^{s-t} E_t[r_{t+s}]$$

- Two value functions:
 1. state-value function: $V(s) = \max_{a \in A} \{E[r|s, a] + \gamma E[V(s')|s, a]\}$
 2. action-value function: $Q(s, a) = E(r|s, a) + \gamma E[\max_{a' \in A} Q(s', a')|s, a]$
- Q-Learning:
 - * Initialize $Q_0(s, a)$. Loop until converge:
 - agent chooses $a = \arg \max Q(s, a)$ and observes $\{r, s', s, a\}$.
 - Old Q: $Q(s, a)$.
 - Target Q: $r_t + \gamma \max_{a' \in A} Q(s', a')$.

- New Q:

$$Q'(s, a) = (1 - \alpha)Q(s, a) + \alpha[r + \gamma \max_{a' \in A} Q(s', a')]$$

- Both in theory and in practice, one of the most important ingredients for successful learning is exploration. Agents do not always act according to what they currently believe to be optimal but they sometimes choose actions at random in order to learn. How much exploration to perform is an open problem, called the exploration-exploitation dilemma.
- Proof of convergence and theory.
 - Tsitsiklis [1994] perform an extensive convergence analysis of q-learning algorithms. His proof of convergence is based on contraction mapping theorem.
 - Bertsekas and Yu [2012] introduce a mix of policy iteration and q-learning that improves efficiency even for approximate methods. It seems great check it.
 - Ma et al. [2020] show formally that the recursive specification of the q values can be considered as a transformation of the bellman operator. The authors use this q-transform to solve problems with unbounded rewards.
 - Ma and Stachurski [2021] use computational complexity theory and numerical experiments to conclude that the q-transform leads to gains in computational efficiency. NOTES: WHY?

2.1 Deep Reinforcement Learning

- Deep Q Network (DQN)
 - DQN consists in approximating the Q function with a neural net: $\hat{Q}(s, a) \equiv \mathcal{N}_\rho \sim Q(s, a)$.
 - The parameters ρ are estimated in order to minimize

$$\mathcal{L}(\rho) = \mathbb{E} \left[\left(\hat{Q}(s, a) - (r_t + \gamma \max_{a' \in A} \hat{Q}(s', a')) \right)^2 \right]$$

in observed transitions $\{r, s', s, a\}$.

- The most widely used NN is a composite function made of nested linear functions:

$$\mathcal{N}_\rho(x) = \sigma_K(W_K \dots \sigma_2(W_2 \sigma_1(W_1 x + b_1) + b_2) \dots + b_K)$$

where x is the input (in our case the transitions $\{r, s', s, a\}$), $W_i \in \mathbb{R}^{m_{i+1} \times m_i}$ are the weights, $b_i \in \mathbb{R}^{m_{i+1}}$ are the biases, and $\sigma()$ are activation functions.

2.2 State of the Art Algorithms

- Explain SAC or PPO.

2.3 Multi-Agent Reinforcement Learning

- Define Partially Observed Markov Games.
- Example 1: Algorithmic Collusion.
- Algorithm: Independent learning.
- Challenges of multi-agent.
- Solutions: QMix and other algorithms.
- Prominent success cases:
 - [Berner et al. \[2019\]](#) achieve super-human performance at the game Dota 2. The game is challenging for multiple reasons. It is multi-agent, stochastic, it has a large action and observation space and it requires long term strategy. They authors used a single neural net to approximate both the value function and the policy function. In particular, they used a single-layer 4096-unit LSTM.
 - [Baker et al. \[2019\]](#) show that agents self-playing hide and seek in a complex environment exhibit emergent curricular learning, in which they progressively discover strategies and then devise counter-strategies. They use PPO with GAE. The distributed computation framework is called rapid.
 - [Vinyals et al. \[2019\]](#) achieve grand-master level at the game StarCraft 2. “Observations of player and opponent units are processed using a self-attention mechanism. To integrate spatial and non-spatial information, we introduce scatter connections. To deal with partial observability, the temporal sequence of observations is processed by a deep long short-term memory (LSTM) system. To manage the structured, combinatorial action space, the agent uses an auto-regressive policy and recurrent pointer network.”. “For every training agent in the league, we run 16,000 concurrent StarCraft II matches and 16 actor tasks (each using a TPU v3 device with eight TPU cores²³) to perform inference. The game instances progress asynchronously on preemptible CPUs (roughly equivalent to 150 processors with 28 physical cores each), but requests for agent steps are batched together dynamically to make efficient use of the TPU. Using TPUs for batched inference provides large efficiency gains over previous work^{14,29}. Actors send sequences of observations, actions, and rewards over the network to a central 128-core TPU learner worker, which updates the parameters of the training agent. The received data are buffered in memory and replayed twice. The learner worker performs large-batch synchronous updates. Each TPU core processes a mini-batch of four sequences, for a total batch size of 512. The learner processes about 50,000 agent steps per second. The actors update their copy of the parameters from the

learner every 10 s. We instantiate 12 separate copies of this actor–learner setup: one main agent, one main exploiter and two league exploiter agents for each StarCraft race. One central coordinator maintains an estimate of the payoff matrix, samples new matches on request, and resets main and league exploiters. Additional evaluator workers (running on the CPU) are used to supplement the payoff estimates. See Extended Data Fig. 6 for an overview of the training setup.”

3 Baseline Framework

- N^h households indexed by i choose how much to consume and how much to invest in order to maximize $\sum_{t=0}^{\infty} \beta^t U(C_{i,t})$.

- The final good is produced using N^c capital goods indexed by j according to technology

$$Y_{i,t} = Z_{i,t}^h \left(\prod_{j=1}^{N^c} K_{i,j,t}^{1/N^c} \right)^\alpha = Z_{i,t}^h \prod_{j=1}^{N^c} K_{i,j,t}^{\alpha/N^c}$$

where total factor productivity $Z_{i,t}^h$ follows the stochastic process $Z_{i,t+1}^h \sim \mathcal{P}(Z_{i,t}^h)$.

- The agent allocates $Y_{i,t}$ between consumption and investment in new capital goods. The evolution of the stock of capital is

$$K_{i,j,t+1} = (1 - \delta)K_{i,j,t} + I_{i,j,t}^h$$

where $I_{i,j,t}^h$ represents investment in capital good j by household i in period t ¹.

- The cost in term of the final good of producing $I_{j,t}^h$ new units of capital good j is $\frac{\phi}{2} \left(I_{j,t}^h \right)^2$.
- In order to frame the problem, we define the saving rate $s_{j,t}$ as the expenditure on new capital goods j in terms of the final good. Thus, consumption is

$$C_{i,t} = \left(1 - \sum_{j=1}^{N^c} s_{i,j,t} \right) Y_{i,t}$$

- We will consider both a competitive market’s formulation and a planner’s formulation of the problem:
 - First, we assume that there is a market for investment goods and we introduce capital good firms that pay the adjustment cost and sell the investment good at price $p_{j,t}^k$. Thus, from the point of view of the household, investment is $I_{i,j,t}^h = s_{i,j,t}/p_{j,t}^k$.

¹The superscript h is added because we will consider both planner and market solutions to the problem, and in the latter it is useful to keep track of the agent that is choosing the variable I .

- Second, in order to formulate the planner's problem, we assume that all the resources that are committed by households for investment on capital good j are actually employed in production, and then total production of capital good j is distributed among households according to their contribution. In this case, total production is defined implicitly by:

$$\sum_{i=1}^{N^h} s_{i,j,t} Y_{i,t} = \frac{\phi}{2} I_{j,t}^2 \quad \text{for } j \in [1, \dots, N^h]$$

Solving for $I_{j,t}$ we get:

$$I_{j,t} = \sqrt{\frac{2}{\phi} \sum_{i=1}^{N^h} s_{i,j,t} Y_{i,t}}$$

We can then distribute it among households according to

$$I_{i,j,t} = \frac{s_{i,j,t} Y_{i,t}}{\sum_{i=1}^{N^h} s_{i,j,t} Y_{i,t}} I_{j,t} = \frac{s_{i,j,t} Y_{i,t}}{\sqrt{\frac{\phi}{2} \sum_{i=1}^{N^h} s_{i,j,t} Y_{i,t}}}$$

3.1 Competitive Market's Formulation

- Now we have two types of agents. The household purchases capital goods and consume the final good. The capital good firm, which produces the capital good subject to a quadratic cost.

3.1.1 Household

- The recursive formulation of the problem for household i is:

$$\begin{aligned} V_i(\{K_{i,j,t}\}_{i,j}) &= \max_{\{s_{i,j,t}, K_{i,j,t+1}\}_j} U(C_{i,t}) + \beta \mathbb{E}_t V_i(\{K_{i,j,t+1}\}_{i,j}) \quad \text{s.t.} \\ C_{i,t} &= (1 - \sum_{j=1}^{N^c} s_{i,j,t}) Y_{i,t} \\ Y_{i,t} &= Z_{i,t}^h \prod_{j=1}^{N^c} K_{i,j,t}^{\alpha/N^c} \\ K_{i,j,t+1} &\leq (1 - \delta) K_{i,j,t} + \frac{s_{i,j,t} Y_{i,t}}{p_{j,t}^k} \quad \text{for } j \in [1, \dots, N^c] \end{aligned}$$

- The lagrangian of the problem is:

$$\mathcal{L}_t = \mathbb{E}_t \sum_{r=0}^{\infty} \beta^r \left(U \left(\left(1 - \sum_{j=1}^{N^c} s_{i,j,t+r} \right) Z_{i,t+r}^h \prod_{j=1}^{N^c} K_{i,j,t+r}^{\alpha/N^c} \right) + \sum_{j=1}^{N^c} Q_{i,j,t+r} \left[(1 - \delta) K_{i,j,t+r} + \frac{s_{i,j,t+r} Z_{i,t+r}^h \prod_{j=1}^{N^c} K_{i,j,t+r}^{\alpha/N^c}}{p_{j,t+r}^k} - K_{i,j,t+r+1} \right] \right)$$

- The first order conditions are:

$$\begin{aligned} [s_{i,j,t}] \quad Q_{i,j,t} &= p_{j,t}^k U'(C_{i,t}) \\ [K_{i,j,t+1}] \quad Q_{i,j,t} &= \beta \mathbb{E}_t \left(U'(C_{i,t+1}) \left(1 - \sum_{j=1}^{N^c} s_{i,j,t+1} \right) \frac{\alpha}{N^c} \frac{Y_{i,t+1}}{K_{i,j,t+1}} \right. \\ &\quad \left. + Q_{i,j,t+1} \left[(1 - \delta) + \frac{s_{i,j,t+1}}{p_{j,t+1}^k} \frac{\alpha}{N^c} \frac{Y_{i,t+1}}{K_{i,j,t+1}} \right] \right) \end{aligned}$$

- Combining both F.O.C.s we get the inverse demand:

$$p_{j,t}^k = \mathbb{E}_t \left(\beta \frac{U'(C_{i,t+1})}{U'(C_{i,t})} \left[\left(1 - \sum_{j=1}^{N^c} s_{i,j,t+1} + s_{i,j,t+1} \right) \frac{\alpha}{N^c} \frac{Y_{i,t+1}}{K_{i,j,t+1}} + p_{j,t+1}^k (1 - \delta) \right] \right) \quad (1)$$

3.1.2 Capital good firms (superscript c)

- Capital good firms maximize the discounted sum of profits. Profits in each period are

$$\pi_{j,t}^c = p_{j,t}^k I_{j,t}^c - \frac{\phi}{2} (I_{j,t}^c)^2$$

- The lagrangian of the problem for firm j is:

$$\mathcal{L}_{j,t} = \mathbb{E}_t \sum_{r=0}^{\infty} \beta_f^r \left(p_{j,t+r}^k I_{j,t+r}^c - \frac{\phi}{2} (I_{j,t+r}^c)^2 \right)$$

- The first order conditions is:

$$[I_{j,t}^c] \quad p_{j,t}^k = \phi I_{j,t}^c \quad (2)$$

3.1.3 Market clearing and equilibrium definition

- Given that $I_{i,j,s}^h = s_{i,j,t} Y_{i,t} / p_{j,t}^k$, we can write the market clearing condition as:

$$\sum_{i=1}^{N^h} \frac{s_{i,j,t} Y_{i,t}}{p_{j,t}^k} = I_{j,t}^c \quad \text{for } j \in [1, \dots, N^c]$$

- **Equilibrium:** A competitive equilibrium is a sequence of prices $\{p_{j,t}^k\}_{j,t}$ and allocations $\{s_{i,j,t}, I_{j,t}^c\}_{i,j,t}$ such that:
 - Given prices, household maximize the intertemporal utility of consumption and capital good firms maximize intertemporal profits.
 - Market for capital goods $j \in [1, \dots, N^c]$ clear.

3.1.4 System of Equations

- The system of equations is:

$$\begin{aligned} p_{j,t}^k &= \mathbb{E}_t \left(\beta \frac{U'(C_{i,t+1})}{U'(C_{i,t})} \left[\left(1 - \sum_{j=1}^{N^c} s_{i,j,t+1} + s_{i,j,t+1} \right) \frac{\alpha}{N^c} \frac{Y_{i,t+1}}{K_{i,j,t+1}} + p_{j,t+1}^k (1 - \delta) \right] \right) \quad \text{for } i, j \\ p_{j,t}^k &= \phi I_{j,t}^c \quad \text{for } j \\ I_{j,t}^c &= \sum_{i=1}^{N^h} \frac{s_{i,j,t} Y_{i,t}}{p_{j,t}^k} \quad \text{for } j \\ K_{i,j,t+1} &= (1 - \delta) K_{i,j,t} + \frac{s_{i,j,t} Y_{i,t}}{p_{j,t}^k} \quad \text{for } i, j \end{aligned}$$

- We want to reduce this system to $N^h * N^c$ equations that only depend on $\{s_{i,j,t}, s_{i,j,t+1}\}_{i,j}$. We start by replacing the market clearing condition in the supply curve:

$$p_{j,t}^k = \sqrt{\phi \sum_{i=1}^{N^h} s_{i,j,t} Y_{i,t}}$$

- We now replace this expression for $p_{j,t}^k$ in the transition law for capital $K_{i,j,t+1}$:

$$K_{i,j,t+1} = (1 - \delta) K_{i,j,t} + \frac{s_{i,j,t} Y_{i,t}}{\sqrt{\phi \sum_{i=1}^{N^h} s_{i,j,t} Y_{i,t}}} \quad (3)$$

(4)

- Finally, we can replace these expressions for $p_{j,t}^k$ and $K_{i,j,t+1}$ in the inverse demand function. Notice

that $Y_{i,t+1}$ depends on the set of $\{K_{i,j,t+1}\}_j$ and $C_{i,t+1}$ is a known function of the set $\{s_{i,j,t+1}, K_{i,j,t+1}\}_j$, so we get $N^c * N^h$ differential equations whose only unknowns are $\{s_{i,j,t}, s_{i,j,t+1}\}_{i,j}$.

3.2 Deterministic Steady state

- In the deterministic steady state, we assume that $Z_i^h = 1, \forall i \in [1, \dots, N^h]$. Then, all households have a symmetric problem and all capital goods are equivalent. Consequently, we can drop subscripts i and j on all variables. The steady state equations become:

$$\begin{aligned} Y &= K^\alpha \\ p^k &= \sqrt{\phi N^h s Y} \\ \delta K &= \frac{sY}{\sqrt{\phi N^h s Y}} \\ p^k &= \frac{\beta}{1 - (1 - \delta)\beta} \left[(1 - (N^c - 1)s) \frac{\alpha}{N^c} K^{\alpha-1} \right] \end{aligned}$$

- The solution is:

$$K = \left[\phi \delta N^h N^c \left(\frac{1 - \beta * (1 - \delta)}{\alpha \beta} + \frac{\delta (N^c - 1)}{N^c} \right) \right]^{\frac{1}{\alpha-2}}$$

- For $N^f = 1, N^c = 1$, The steady state is:

$$K = \left(\frac{\beta_f}{1 - (1 - \delta)\beta_f} \frac{\alpha^f}{\phi \delta} \right)^{\frac{1}{2-\alpha^f}}$$

3.3 Planner's Formulation

- The recursive formulation of the problem is:

$$\begin{aligned} V(\{K_{i,j,t}\}_{i,j}) &= \max_{\{s_{i,j,t}, K_{i,j,t+1}\}_{i,j}} \sum_{i=1}^{N^h} U(C_{i,t}) + \beta \mathbb{E}_t V(\{K_{i,j,t+1}\}_{i,j}) \quad \text{s.t.} \\ C_{i,t} &= (1 - \sum_{j=1}^{N^c} s_{i,j,t}) Y_{i,t} \quad \text{for } i \in [1, \dots, N^h] \\ Y_{i,t} &= Z_{i,t}^h \Pi_{j=1}^{N^c} K_{i,j,t}^{\alpha/N^c} \quad \text{for } i \in [1, \dots, N^h] \\ K_{i,j,t+1} &\leq (1 - \delta) K_{i,j,t} + \sqrt{\frac{2}{\phi}} s_{i,j,t} Y_{i,t} \quad \text{for } i \in [1, \dots, N^h] \text{ and } j \in [1, \dots, N^c] \end{aligned}$$

- The lagrangian of the problem is:

$$\mathcal{L}_t = \mathbb{E}_t \sum_{r=0}^{\infty} \beta^r \left(\sum_{i=1}^{N^h} U \left(\left(1 - \sum_{j=1}^{N^c} s_{i,j,t+r} \right) Z_{i,t+r}^h \prod_{j=1}^{N^c} K_{i,j,t+r}^{\alpha/N^c} \right) + \sum_{i=1}^{N^h} \sum_{j=1}^{N^c} Q_{i,j,t+r} \left[(1 - \delta) K_{i,j,t+r} + \sqrt{\frac{2}{\phi} s_{i,j,t+r} Z_{i,t+r}^h \prod_{j=1}^{N^c} K_{i,j,t+r}^{\alpha/N^c}} - K_{i,j,t+r+1} \right] \right)$$

- The first order conditions are:

$$\begin{aligned} [s_{i,j,t}] \quad Q_{i,j,t} &= \phi \sqrt{\frac{2}{\phi} s_{i,j,t} Y_{i,t}} \frac{\partial U(C_{i,t})}{\partial C_{i,t}} = \phi I_t U'(C_{i,t}) \\ [K_{i,j,t+1}] \quad Q_{i,j,t} &= \beta \mathbb{E}_t \left[\frac{\partial U(C_{i,t})}{\partial C_{i,t}} \left(1 - \sum_{j=1}^{N^c} s_{i,j,t+1} \right) \frac{\alpha}{N^c} \frac{Y_{i,t+1}}{K_{i,j,t+1}} \right. \\ &\quad \left. + Q_{i,j,t+1} \left((1 - \delta) + \left(\frac{2}{\phi} s_{i,j,t+1} Y_{i,t+1} \right)^{-1/2} \frac{s_{i,j,t+1}}{\phi} \frac{\alpha}{N^c} \frac{Y_{i,t+1}}{K_{i,j,t+1}} \right) \right] \end{aligned}$$

To do:

- Change problem to overall cost. b
- Simplify FOC for K . you can use x/\sqrt{x} on the right side and write in terms of investment.
- Calculate steady state system of equations.
- Order the rest of the doc.
- Why do we need the planner for? Just as an analog for the program?

4 Applying Reinforcement Learning

4.1 Learning by a Planner

- Exercise 1: 1 capital good, many households (1 vs 4 vs 100)
- Exercise 2: many capital goods, (1v1, 2v2, 10v10)

4.2 Multi-agent learning

- Explain centralized critic algorithm.
- Explain parallelized computation framework.
- Exercise 1: 1 vs 5 vs 100 households.

- Exercise 2: Capital Good firms as AI agents. 1 vs 1, 2 vs 2, 10 vs 10. Does it converge to competitive equilibrium?
- Exercise 3: Half the households have full info and half the households only observe own and aggregate stock.

5 Design Lessons

- Exercise 1: Choosing spending vs choosing quantities. General normalization issues.
- Exercise 2: How to specify markets.
- Exercise 3: How to impose constraints.
- Exercise 3: Learning rates and algorithmic issues?

6 Conclusion

to be done

References

- Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor Mordatch. Emergent tool use from multi-agent autocurricula. *arXiv preprint arXiv:1909.07528*, 2019.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Dimitri P Bertsekas and Huizhen Yu. Q-learning and enhanced policy iteration in discounted dynamic programming. *Mathematics of Operations Research*, 37(1):66–94, 2012.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International Conference on Machine Learning*, pages 1407–1416. PMLR, 2018.
- Qingyin Ma and John Stachurski. Dynamic programming deconstructed: Transformations of the bellman equation and computational efficiency. *Operations Research*, 2021.
- Qingyin Ma, John Stachurski, and Alexis Akira Toda. Unbounded dynamic programming via the q-learning transform. *arXiv preprint arXiv:2012.00219*, 2020.

John N Tsitsiklis. Asynchronous stochastic approximation and q-learning. *Machine learning*, 16(3):185–202, 1994.

Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.