

ALMA MATER STUDIORUM · UNIVERSITÀ DI BOLOGNA

DEPARTMENT OF ECONOMICS
Second-cycle · Master's Degree
in
ECONOMICS

**A STUDY ON TAXATION POLICIES
ON HETEROGENEOUS
AI-DRIVEN AGENTS**

Defended by:
LORENZO MEZZINI
926012

Supervisor:
VINCENZO DENICOLÒ

Graduation Session of December
Academic Year 2020/2021

Abstract

Throughout the thesis, it will be presented a partial paper reproduction, followed by an extension. This work presents a set of economical simulations with four heterogeneous AI-driven agents. This is done by the use of a python package called ai-economist[1] and techniques from reinforcement learning. The training of the agents is going to be divided into two steps, a first phase where the agents get accustomed to the environment, and a second one where different kinds of taxation will be imposed. The main focus will be the relationship between total production and equality. Quite reasonably, we found that there is a negative impact from taxes on productivity, and a positive impact on equality.

Contents

Abstract	i
Introduction	1
1 What is Foundation	3
1.1 Gather and trade	3
1.1.1 World and entities	4
1.1.2 Game dynamics	5
1.1.3 Agents	6
1.1.4 Policymaker	8
2 Reinforcement Learning	10
2.1 Finite Markov Decision process (MDPs)	11
2.2 Approximate solution Methods	14
2.2.1 Policy gradient Methods	14
2.2.2 Proximal Policy Optimization Algorithms	14
3 Experiments	17
3.1 Experiment Setup	17
3.2 Phase 1 training	19
3.3 Phase 2 Training	22
3.3.1 Free Market	22
3.3.2 US taxation	22
3.3.3 Italian Taxation	23

3.3.4 Communism	23
3.4 Training Results	24
Conclusions	27
A Tables, Algorithms and Graphs	30

List of Figures

1.1	Rendering of the world at initial conditions:	5
2.1	Agent environment interaction dynamic:	11
2.2	One step L^{CLIP} as a function of $r_t(\theta)$:	15
3.1	Learning rate and Labor weighed cost for Phase 1 training: . .	19
3.2	Equality and Productivity during Phase 1:	20
3.3	Phase 1 full brake down:	21
3.4	Phase 2 training brake down:	25
3.5	Productivity and Equality histogram:	26
3.6	Coin Distribution	26
3.7	Coin Distribution	26
3.8	Coin equality times total production	27
A.1	Free market full brake down:	35
A.2	US full brake down:	36
A.3	Italy full brake down:	37

A.4 Communism full brake down: 38

List of Tables

3.1 Fully connected neural network utilized. 18

3.2 2018 US federal tax system: 23

3.3 2020 INPS tax system: 23

A.1 Hyperparameters for the environment. 31

A.2 Hyperparameters for the Training. 32

A.3 Full observation space. 33

List of Algorithms

1 PPO with clipped objective 16

2 Agents Learning Loop. 34

Introduction

One of the challenges in macroeconomic theory is to transpose the theory to the real world. In particular the process of designing a policy and implementing it. This is because the analytical solutions are bound to certain a level of complexity. The more complex and similar to reality the problem gets the harder it gets to solve it analytically. In addition agents in the real world are heterogeneous and thus each one is acting differently according to his skill and initial conditions.

In this thesis, I am going to present a partial reproduction with an expansion of the paper *The AI Economist: Optimal Economic Policy Design via Two-level Deep Reinforcement Learning*[2]. With the scope of presenting an AI-driven economical simulation, where is possible to investigate a high complexity setup that would be almost impossible to study analytically. This simulation requires the agents to behave in a complex way. They are asked to maximize their utility in a simulated economy where they can gather goods, use them or interact with other agents in a bid/ask market. To achieve this goal we will use a technique called proximal policy optimization. This is an optimization algorithm that will train a fully connected neural network. It will be shown that AI-driven agents can find a sub-optimal behavior, and come up with emerging behaviors and specialization to maximize their utility.

Afterward, multiple kinds of taxation will be introduced in the simulations as a discriminant factor, and we will compare the results on total economical productivity and coin equality among agents. This requires a second training

that will get the agents accustomed to the new setup.

The paper that is used as a reference [2], addresses the question of whether is it possible to generate optimal policies through the use of reinforcement learning. To do it, they first built a gather-and-trade game, publicly available on GitHub, as a framework for the simulations. Afterward, they trained the agents to behave optimally within this simulation. Once they had this setup they created four simulations: one that recreates the US taxes, the free market scenario, the Saez taxation, and another RL driven policymaker that optimize the production of equality and productivity. They were able to show that the RL agent produced an economy the 16% more efficient compared with the analytical solution purposed by Saez[3], in the variable of interest. The entirety of the code from this second part is not available since their paper is still in peer review.

Thus, using third-party code available on Github, I will try to reproduce the first step of optimization, where the agents maximise their utility. Then I will introduce four taxations as well, however, two of these will be the US and the free market, and the other two will be the Italian system and communism. Hence, the main questions that this thesis address are: Is reinforcement learning a viable media to construct an economical accurate simulation? Can we conclude the impact of the different taxation systems on the overall economy?

Chapter 1

What is Foundation

In this chapter, I am going to introduce the package developed by Alex Trott and Stephen Zheng and provided by Salesforce called Foundation[4]. This package offers the possibility to create simulations with multiple agents that interact in a 2D world. At first, there will be a presentation of the gather-and-build simulation with the relevant entities and dynamics. Furthermore, the basic agents will be dismantled in their main components to better understand their action space and scope, and finally, there will be a small discussion on the behavior of the policymaker.

1.1 Gather and trade

The simulation that will be used through the dissertation is called gather and trade. This simulation takes place in a 2D map that represents the world where the agents live and interact. The shape of the world is a 25x25 grid where are disposed various kinds of entities. Within this world, 4 agents are free to move around, gather resources, trade them and use them to build houses. These agents are different for their skill level, allowing them to have higher/lower rewards for their actions. A fifth agent, called policymaker, is asked to tax and redistribute the income of the four previous agents, based on information

about their endowments, but not their skill.

1.1.1 World and entities

As aforementioned the map is a 25x25 grid where are present some entities. Some of these are visible in the world as water or stone, others are just present in the agent's endowment as coins or labour. The complete set of entities is:

- Water
- House
- Source Block (Wood and Stone)
- Wood and Stone
- Coins
- Labor

Water is a simple world block that has no use other than avoiding agents passing through. We can see from Figure 1.1 that the water is used to divide the world into 4 macro-areas, each one with different resources. A House is a block that is not present at the beginning of the simulation, but it can be built by agents in an empty block, only the owner of the house can walk through it.

Source blocks are two entities that spawn stochastically resources, namely wood, and stone, as we can see from Figure 1.1 in the four areas divided by the water we have a zone with both wood and stone source block, two other areas with just one kind of source block and one that has no resources at all.

Coins are the measurement of the value produced in the world. Coins are generated when a house is built, the agent that builds the house is rewarded with a certain amount of coins that varies with the skill level.

Labor is a measurement of the total effort exerted by an agent, it is generated every time an agent takes an action and is perceived as a dis-utility from the agent.

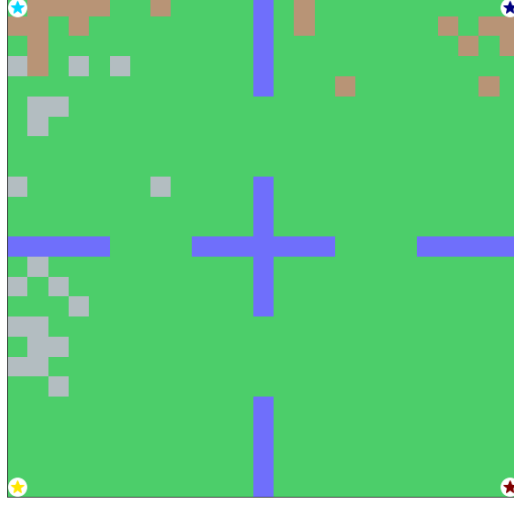


Figure 1.1: Rendering of the world at initial conditions: *this is a rendering of the world map at the first timestep of the simulation, blue blocks are water, brown blocks are wood sources, grey blocks are Stone sources. The four corners are the starting positions for the 4 agents.*

1.1.2 Game dynamics

A problem that has a continuous flow of agent-environment interactions can be formalized by a Finite Markov Decision problem [5]. In particular, this simulation is a partial-observable multi-agent Markov Games (MGs). The problem is defined by the tuple $(S, A, R, \gamma, o, \mathcal{T})$ where S is the state space, A is the action space, R the reward space and gamma a dicount factor.

The simulation is composed of a series of episodes, each of length H time-steps. At every point in time $t \in [0, H]$ the episode is characterized by a state s_t , representing the current world environment, every agent performs an action $a_{i,t}$ given the current partial observation $o_{i,t}$ of the world state, and receives a reward r_{it} . Afterwards the environment transits to a new state s_{t+1} according to the transition distribution $\mathcal{T}(s_{t+1}|s_t, \mathbf{a}_t)$. This chain of interactions state-action-reward/state carries on until the end of an episode.

$$s_0 \rightarrow_{o_0} \mathbf{a}_0 \rightarrow r_1, s_1 \rightarrow_{o_1} \mathbf{a}_1 \rightarrow \dots \rightarrow_{o_{H-1}} \mathbf{a}_{H-1} \rightarrow r_H, s_H$$

Here \mathbf{a} and \mathbf{o} are the vectorized observations and actions for all the five agents. Given the particular structure of the simulation every single agent will receive an observation at every time-step (different for everyone, more on that later), but only at the 4 basic agents will be asked to act, the policymaker will act only upon a certain condition met. In this case, the episode lasts for 1000 time-step and the policymaker is asked to act (tax the other agents) every other 100 steps. The existence of multiple episodes is necessary for the 4 agents and the policymaker to define their own optimal policy $\pi_i(o_{i,t}, h_{i,t-1}; \theta_i^*)$, this optimization process will be discussed in the RL chapter.

1.1.3 Agents

From the information above we know that the four basic agents are endowed with labor, coins, wood and stone. They live in the world map, can act within it and their objective is to maximize their γ -discounted utility function. Now I will describe in more detail the agents starting from the information that they receive at each time-step, then talking about the actions that they are allowed to take and finally about their objective.

Observation space: Given that this simulation is a partial-observable multi-agent Markov Game, the observation that agent i receive at time t is not complete but partial, this information can be summarized in the following way:

- $o_{i,t}^{\text{world state}}$: world map situation surrounding the agent, this is limited to the 11x11 grid around the agent i .
- $o_{i,t}^{\text{market state}}$: full information about the market state for wood, stone and available trades.
- $o_{i,t}^{\text{agent}}$: public resources and coin endowments (this information is also available to the policy maker), private labor performed and skill level.
- $o_{i,t}^{\text{tax}}$: tax structure

- $o_{i,t}^{\text{other}}$: other information (ex. action mask)

the full observation space can be seen in Table A.3

Action space: The agent can take one action per time-step and can choose this action from the 50 listed below:

- **Movement:** 4 actions for the basic movements N, S, E, W
- **Gather:** 1 action for gathering resources
- **Trade:** 44 actions for trading resources
- **Idle:** 1 action that does nothing

The movement's actions along with gather do not need much of an explanation, these are simple actions that increases labor by 0.21 units each time pursued. The building action requires the agent to consume (destroy) one unit of wood and one unit of Stone, as a consequence he gains 2.1 units of labor and an amount of coin that depends on his skill level. The most complicated set of actions are the one that rules trading. Each one of them is a combination of the 11 price levels [0,1,...,10] that the agent is willing to (pay/request) for each side (bid/ask) and for each resource (wood/stone). A trading open action remains open for 50 turns, if in this time span it is matched by the corresponding counteraction at the same price (a bid for an ask and vice versa) then the trade takes place and each agent gets 0.05 unit of labor.

Furthermore, notice that during the episode agents might incur into a series of inconclusive actions, such as moving north while at the north border of the map or building a house without the required wood and stone. The so called action mask, present in the observation space, is used in the learning process to avoid wasting time exploring these possibilities. It is an array of binary values of length 50 that "masks" all meaningless actions.

Agent objective Agents in the simulation earn coins when building houses or trading goods, the utility for the four agents is an isoelastic utility:

$$u_i(x_{i,t}, l_{i,t}) = crra(x_{i,t}^c) - \vartheta_k l_{i,t}, \quad crra(z) = \frac{z^{1-\eta} - 1}{1 - \eta}, \quad \eta > 0 \quad (1.1.1)$$

Where $l_{i,t}$ is the cumulative labor associated with the actions taken up to time t , $x_{i,t}^c$ is the coin endowment and ϑ is a function utilized for labor annihilation with the following specification $\vartheta_k = 1 - \exp\left(-\frac{\text{number of time mean reward} \geq 0}{\text{energy warm-up constant } (k)}\right)$. This variable will play an important role during the two-step optimization process purposed in the original paper. In particular, during phase 1 of training, the labor cost is annihilated to help agents avoid sub-optimal behaviors. And η determines the degree of non-linearity of the utility. This utility function is assumed to be the same for all the agents.

The maximization problem is solved for a rational behaving agent by optimizing the total discounted utility over time,

$$\forall i : \max_{\pi_i} \mathbb{E}_{a_i \sim \pi_i, \mathbf{a}_{-i} \sim \pi_{-i}, s' \sim \mathcal{T}} \left[\sum_{t=1}^H \gamma^t r_{i,t} + u_i(x_{i,0}, l_{i,0}) \right] \quad (1.1.2)$$

with $r_{i,t} = u_i(x_{i,t}, l_{i,t}) - u_i(x_{i,t-1}, l_{i,t-1})$ being the instantaneous reward of agent i at time t . Equation 1.1.2 illustrates a multi-agent optimization problem in which actors optimize their behavior at the same time since the utility of each agent is dependent on the behavior of other agents. Another agent, for example, may deny an agent access to resources, limiting how many houses the agent can create in the future and hence its utility. While computing equilibrium for complicated environments like this is still out of reach, we will see later how RL may be utilized to produce meaningful, emergent behaviors.

1.1.4 Policymaker

The policymaker, or social planner, differs deeply from the previous agents. Being the focus of the research question its structure and behavior change a lot in every single simulation. Due to the complexity of creating a multi-agent neural network with different objectives I was not able to introduce

a policymaker that behaves according to an RL algorithm. This makes the policymaker less interesting from a structural point of view but its behavior is still meaningful for the results it produces.

The only observation needed is the number of coins that each agent possesses. With this value is possible to place an agent within a tax bracket and tax it accordingly to some predefined values.

There will be four simulations in this thesis:

- US taxation
- Italian taxation
- Free market
- Communism

and the policymaker will act accordingly to these taxation systems. The US taxation is based on the 2018 US federal tax and brackets, and the Italian is based on the 2020 Irpef tax and brackets. The brackets are created in such a way that 1 coin in the simulation is equal to 1000 \$ or €863.93 according to the tax system.

Chapter 2

Reinforcement Learning

Reinforcement learning is the process of learning what to do to maximize a numerical reward. A mechanism where the learner is not told what action to take, but instead must discover which action yields the highest reward by trying them. In this case, actions might affect not only the immediate reward but also future situations and all the subsequent rewards. These two characteristics – Trial-and-error search and delayed reward – are the two most important distinguishing features of reinforcement learning.

Reinforcement learning differs from supervised learning since training is not based on an external dataset of labeled examples, where each situation (observation) is labeled with the correct action to perform (often identify a category). RL, although one might erroneously think the opposite, is also different from unsupervised learning. The main objective for unsupervised learning is to find hidden structures in an unlabeled dataset, whereas RL's main objective is to maximize a reward signal.

In this chapter, we will see the definition of a Markov decision process and those of the state and action value function. Afterward, we will observe how optimization is carried out in an approximate solution method, in particular in the case of the proximal policy optimization since it is the technique utilized

for training the neural network.

2.1 Finite Markov Decision process (MDPs)

Finite Markov decision processes are a class of problems that formalize subsequent decision making, where not only the influence of the action is exerted on immediate reward but also on those in the future. MDP's are suited for RL since they frame the process of learning through repeated interaction between an agent (decision maker), and an environment (ruled by fixed state transition function).

More specifically an agent is asked to take an action $a_t \in \mathcal{A}$ at every time step $t = 0, 1, \dots$. To do so the agent is provided with an observation of the current environment's state $s_t \in \mathcal{S}$ and a reward $r_t \in \mathcal{R}$ from the previously performed action. Afterwards the environment update its state following a transition distribution $\mathcal{T}(s_{t+1}|a_t, s_t)$ and a numerical reward $r_{t+1} \in \mathcal{R} \subset \mathbb{R}$. This process is reproduced every subsequent timestep, this concatenation of interaction is a MDP.

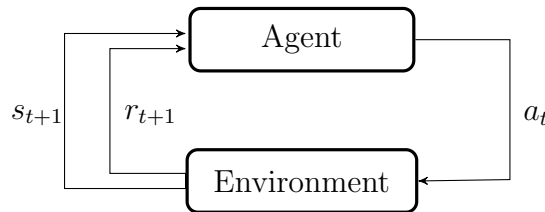


Figure 2.1: Agent environment interaction dynamic: *this is the basic dynamic of the interaction between the agent and the environment in a MDP, at time t the agent provides an action a and the environment responds with a state s and reward r that are used from the agent to decide the action at $t+1$ and so on.*

Objective and Rewards: The main objective of RL is to maximize the total number of rewards it receives. This reward r_t is passed from the environment to the agent at every timestep as a consequence of his actions. In the case of the Gather and Trade the reward is $r_{i,t} = u_i(x_{i,t}, l_{i,t}) - u_i(x_{i,t-1}, l_{i,t-1})$.

Since the agents want to maximize the total upcoming reward we can write this value simply as the sum of all the future rewards.

$$U_t \doteq r_{t+1} + r_{t+2} + \dots + r_H,$$

Where H is the total length of the episode, and at the time step $t = H$ the episode ends. This state is called the terminal state and is a particular state because regardless of the final condition of the agent it reset the environment to the initial condition and restart completely the episode. Another specification for the total reward can implement the concept of discounting, which is more appropriate in the case of economical simulations, thus within the experiments the agent has to maximize his future discounted utility:

$$U_t \doteq r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^{H-t-1} r_H, \quad (2.1.1)$$

the discount rate $\gamma \in [0, 1]$ determines the value that the agent assigns to the future reward, a reward received at k timesteps in the future is only valued γ^{k-1} times what it would be valued today. When the value of γ approaches 0 the agent is more "myopic" and puts most of his interest in immediate rewards, while if it approaches 1 the interest is more projected in the future due to the stronger impact of future rewards.

Value functions: One of the most important elements of RL is the value function. The estimation of this function is one of the crucial points of RL, in fact it tries to quantify the expected return of the rewards it expects to receive. Furthermore, the expected reward depends on the action that the agent decides to take, thus the value function are defined in terms of policies, which are acting behaviors.

If the agent is following policy π at time t , then $\pi(a|s)$ is the probability that the agent takes the action $a_t = a$ given the state $s_t = s$. The aim of RL is

to change the policy based on experience across episodes to find an optimal behavior.

We can write a value function for a state s under the policy π . This function is the expected return when the initial state is s and the policy followed is π .

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi} [U_t | s_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^H \gamma^k r_{t+k+1} \middle| s_t = s \right], \quad \text{for all } s \in \mathcal{S} \quad (2.1.2)$$

$v_{\pi}(s)$ is called the state-value function for policy π . Following from this equation, it is possible to define the value of taking an action a in the state s following the policy π :

$$q_{\pi}(s, a) \doteq \mathbb{E}_{\pi} [U_t | s_t = s, a_t = a] = \mathbb{E}_{\pi} \left[\sum_{k=0}^H \gamma^k r_{t+k+1} \middle| s_t = s, a_t = a \right], \quad (2.1.3)$$

and $q_{\pi}(s, a)$ is called the action-value function for policy π . The important concept here is that the value functions in 2.1.2 and 2.1.3 can be estimated from experience.

There are multiple strategies to determine an optimal behavior starting from the evaluation of these functions, we can divide these ways in two main groups, tabular solution methods and approximate solution methods[5]. For the former we have Montecarlo methods, Dynamic programming, Temporal-difference learning, n-step bootstrap and others. While for the latter we have On/Off-policy methods with approximation and policy gradient methods.

For the purpose of this thesis we are going to focus only on policy gradient methods and a particular set of optimization policy called proximal policy optimization.

2.2 Approximate solution Methods

The approximate solution methods are sets of strategies thought for those problems, such as ours, where the set of possible states is enormous. It is very likely that every state encountered in a simulation will never have been encountered before. Thus, to make meaningful decisions there is the need to be able to generalize from previous state that are, to some extent, similar. This is accomplished by borrowing the existing generalization theory, usually by using the process of function approximation. However, we are going to use policy gradient methods that do not necessarily need to estimate a value function through function approximation.

2.2.1 Policy gradient Methods

Policy gradient methods are a set of parameterized policies that can select actions without the use of a value function. The method consists in calculating the estimator of the policy gradient and feeding it into a stochastic gradient ascent algorithm. We denote with $\pi_\theta(a|s)$ the stochastic policy such that

$$\pi_\theta(a|s) = \pi(a|s, \theta) = Pr \{a_t = a | s_t = s, \theta_t = \theta\} \quad (2.2.1)$$

for the probability of choosing action a at time t given the state s_t and the parameters θ . The aim of the policy gradient methods is to maximize a performance measure $J_t(\theta)$. This is done by updating the parameters θ with a gradient of the performance itself:

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)} \quad (2.2.2)$$

2.2.2 Proximal Policy Optimization Algorithms

Proximal policy optimization algorithms (PPOs) are a family of policy gradient methods for reinforcement learning, which alternate between sampling data through interaction with the environment, and optimizing a "surrogate" objective function using stochastic gradient ascent.

PPO utilizes as a performance measure the advantage $A(s, a) = q(s, a) - v(s)$, which represents how good an action is compared to the average action in a specific state.

It also relies on the ratio of the policy that we want to refine to the older policy:

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \quad (2.2.3)$$

The surrogate objective function that we optimize with PPO is

$$\mathcal{L}^{CLIP}(\theta) = \hat{E}_t \left[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right] \quad (2.2.4)$$

The maximization of this surrogate function takes the advantages of trust region policy optimization (TRPO) [6], adjusting the parameters in the direction of the ratio $r_t(\theta)$, and maintaining this change bounded in a "clipped" region to avoid drastic refinements of the policy. This is shown well in the Figure 2.2

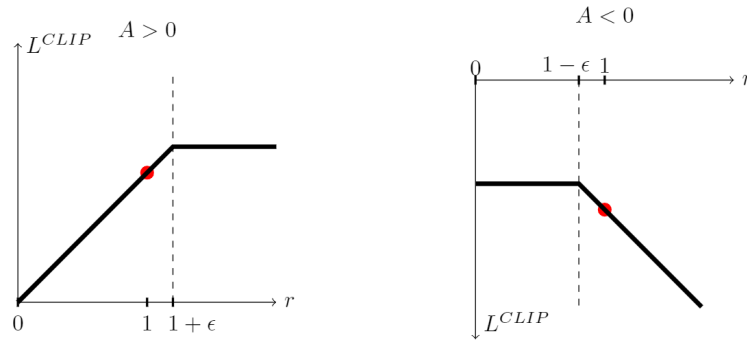


Figure 2.2: One step L^{CLIP} as a function of $r_t(\theta)$: this is the function L^{CLIP} as a function of the probability ratio $r_t(\theta)$ when the advantage A is positive or negative. The red dot is the starting point of the optimization. (Figure from [7])

The surrogate function can be further augmented by adding an entropy bonus to ensure sufficient exploration. By adding these terms, we obtain the following objective function:

$$L_t^{CLIP+VF+S}(\theta) = \hat{E}_t [L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_\theta](s_t)] \quad (2.2.5)$$

where S is an entropy bonus, L_t^{VF} is a squared-error loss and c_1, c_2 are coefficients.

The pseudocode below better explains all the procedure done by the PPO in refining the parameters.

Algorithm 1 PPO with clipped objective

Input: initial policy θ_0 , clipping threshold ϵ

for $k = 0, 1, 2 \dots$ **do**

 Collect a set of trajectories \mathcal{D}_k on policy $\pi_k = \pi(\theta_k)$

 Estimate the advantages $\hat{A}_t^{\pi_k}$, using any advantage est. alg.

 ▷ we will use GAE[6]

 Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$$

by taking K steps of minibatch SGD (via Adam [8]), where

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = E_t \left[\sum_{t=0}^T \left[\min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

end for

Chapter 3

Experiments

In this chapter, I will present the two phases learning process utilized by a fully connected neural network (FCNet) to maximize the agents' utility. During phase 1, the agents' objective will be to maximise utility when labour is gradually introduced. While in phase 2, there will be the introduction of a tax. In this chapter, I will present four taxation mechanisms, two of which are the same as the paper of reference, while the others are an extension of mine. In conclusion, I will present the results derived from the training and a discussion on the values obtained.

3.1 Experiment Setup

After understanding the Foundation framework and the process used by RL to optimize a policy, we can focus on the model used. In the reference paper, the authors used a combination of a convolutional neural network and an LSTM [4]. This could grant them good spatial information from the CNN and memory cells from the LSTM. These features are appropriate because there is the need of processing a map and the agents share the same network.

In the experiments, however, it was used a fully connected neural network that does not share the same features as the one above. The complete structure of

the model is shown in Table 3.1.

Model: "FCNet"			
Layer (type)	Output Shape	Param #	Connected to
observations (InputLayer)	[(None, 1260)]	0	
encoder (Encoder)	(None, 256)	389632	observations[0][0]
fc_out (Dense)	(None, 50)	12850	encoder[0][0]
value_out (Dense)	(None, 1)	257	encoder[0][0]
Total params: 402,739			
Trainable params: 402,739			
Non-trainable params: 0			

Table 3.1: Fully connected neural network utilized.

This setup might incur a loss of efficiency for plenty of reasons. First of all, as said before, the features of the LSTM/CNN are missing. Following, due to lack of resources, it was not possible to do a proper hyperparameter tuning. However, as we will see, the FCNet can provide us with meaningful results, in particular, the agents will show emerging behaviours and specialization. These are characteristics that were also found in the paper by Zheng and Trott, the difference stands in the agents' efficiency which is lower for the model used here.

There are many parameters that have to be set at this stage (for a semi-complete list check the Appendix Table A.1 and Table A.2). The first part of such values is related to the training, while the second part is to the environment. I would like to stress some of the latter parameters, to better understand the simulation.

First, the skill levels are set to be different for each agent. The skill is the number of coins received when building a house. The four skills are fixed and they are taken from a Pareto Distribution. Namely the values are (11.3; 13.3; 16.5; 22.2). These values are always assigned to the same starting location in the map, thus the agent starting in the bottom right will always have the

highest skill, and so on. Other important parameters are the episode length H which is set to 1000 time steps, the resource re-spawn probability which is set to 1% per time step, and the initial coin endowment that is set to 0.

Once the environment is defined, it is possible to start the training which will occur in two phases. The first phase introduces labour, while the second introduces taxation.

3.2 Phase 1 training

The first training phase is necessary to get the FCNet in Table 3.1 accustomed to the world dynamics and avoid falling into un-optimal behaviour once disincentives factors are introduced. These two factors, as already discussed, are labour costs and taxation. At this moment, I am interested in getting the agent accustomed to the labour costs, while the taxes are completely removed and will be addressed in the second phase of the training.

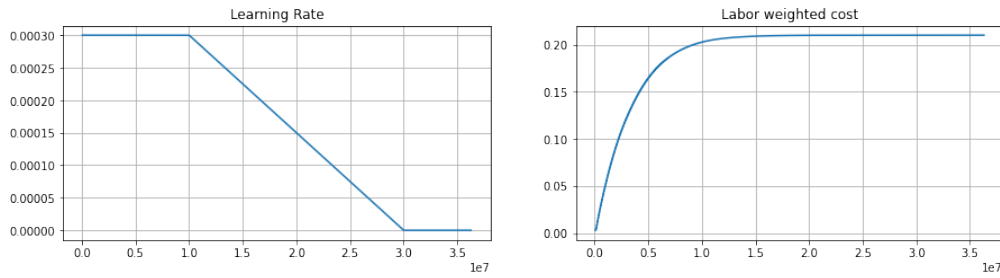


Figure 3.1: Learning rate and Labor weighed cost for Phase 1 training: *this was the schedule for respectively the learning rate and the labor cost in the phase 1 of the training, here the learning rate goes to 0 because we don't seek convergence now.*

The training agenda is built to run for a total of 30M time-steps, as we can see from Figure 3.1 the learning rate is set to be at $3e-4$ for 10M steps, then it linearly reduces to 0 in 20M time-steps. Meanwhile, the labour weight increases following the function:

$$\vartheta_k = 1 - \exp\left(-\frac{\text{number of time mean reward} > 0}{\text{energy warm-up constant } (k)}\right) \quad (3.2.1)$$

where the warm-up constant is set to 10000. This setup gives the FCNet time to learn how to respond properly to the dis-utility generated by the labour.

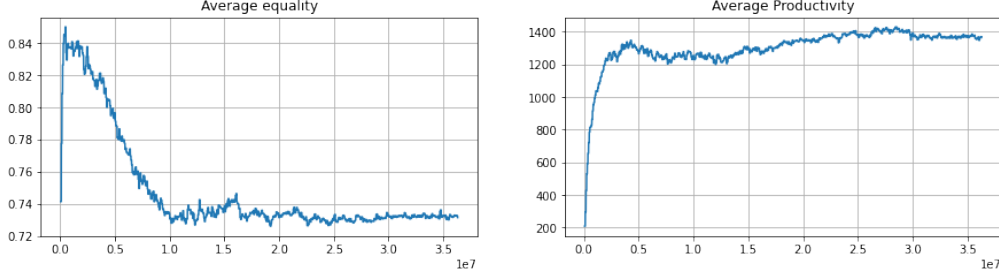


Figure 3.2: Equality and Productivity during Phase 1: *here we have the evolution of equality between agents and total production in the training of phase 1*

In Figure 3.2, we can observe two of the most important variables we are considering: equality and productivity. On the left, we have the average equality, calculated by the following formula:

$$\text{equality}(\mathbf{x}_t^c) = 1 - \frac{N}{N-1} \frac{\sum_{i=1}^N \sum_{j=1}^N |x_{i,t}^c - x_{j,t}^c|}{2N \sum_{i=1}^N x_{i,t}^c} \quad (3.2.2)$$

with $0 < \text{equality}(\mathbf{x}_t^c) < 1$. This function returns 1 if the endowments are equally split between the N (4) agents, and 0 if one agent owns all of the coins in the economy. What we observe is that after 36M time-steps the equality settles around 0.73, this is justified by the difference in coin endowment between the agent with the highest skill (Agent 0 in Figure A.4) and the other agents. Furthermore, we can notice that the average productivity of the economy settles around 1400. However, both of these values are incorrect because of a data collection issue in the code. We will discuss later the actual values.

Regardless, the results obtained so far are not important for the final analysis because there was no criterion used to stop the training and the learning rate

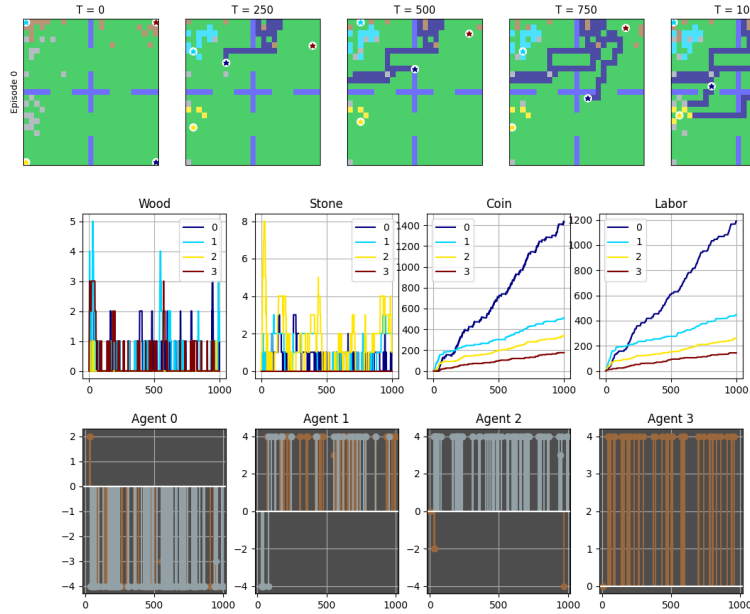


Figure 3.3: Phase 1 full brake down: *this is a more complete overview of one random simulation crated using the weights at the end of phase 1 training*

was decreased to 0. This means that the FCNet might not have reached convergence nor an optimal policy. Nonetheless, we can keep these results in mind as a benchmark for the phase 2 training, as we will use the weights obtained to carry on the training.

In conclusion, the most important result from the first training is that agents show emerging behaviours. We can see that the agents specialize in a specific task according to their initial position on the map and their skill level. Behaviours are important to the scope of the thesis because it is hard to capture specialization analytically. Here, instead, it is rather easy to have specialized agents as their utility is maximized through experience.

From Figure A.4, we can see that the agent with the highest skill has specialized in building houses, this makes sense since he earns more from the building action. By building he can also afford to buy goods from the market, giving other agents the possibility to specialize in gathering and selling. Indeed, this behaviour is adopted by agents 2 and 3, as they are those endowed with the

lowest skills. On the other hand, agent 1 exhibits a variable behaviour. He sells goods in the market, given that he has both of them available in his macro-area, but at the same time he exerts some labour to build houses.

3.3 Phase 2 Training

In the second phase of the training, the learning rate schedule is changed to be $3e-4$ for the first 35M time-steps and then will linearly decrease to $1e-6$ in a span of 15M time-steps. This will grant a faster learning rate in the initial part of the training when the FCNet will suffer a huge decrease in efficiency due to the introduction of the taxation.

The starting point for this second phase is the weights obtained from phase 1, and the training will be formed by four different simulations diverging from this initial state. These simulations are US taxation, Italian Taxation, free-market and Communism.

3.3.1 Free Market

Free-market is the first and simplest kind of simulation that we can carry out starting from the result in phase 1. This training consists in keep training the FCNet without imposing any kind of taxation.

It is reasonable to believe that the result of this first training should generate a more productive economy. Or at least is what should be expected from economical theory [9]. However, the equality between agents might suffer from higher production.

3.3.2 US taxation

The second kind of taxation used is US taxation. This is based on the 2018 Federal brackets and percentages. In particular we can see from Table 3.2 the detailed structure.

Bracket in \$	0-9700	9700-39475	39475-84200	84200-160725	160725-204100	204100-510300	510300+
Bracket in Coin	0-9.7	9.7-39.5	39.5-84.2	84.2-160.7	160.7-204.1	204.1-510.3	510.3+
Tax*	10%	12%	22%	24%	32%	35%	37%

Table 3.2: 2018 US federal tax system: *this is the proportional system that was in place in the US in 2018, notice that the tax is marginal (i.e. if you are in the second bracket you will pay 970\$ + 12% of the amount above 9700\$)*

With the introduction of these lump-sum tax plus redistribution is easy to conclude that there might be an increase in equality, however, this might produce a negative effect on total production since the high-skilled agent are disincentivised to produce.

3.3.3 Italian Taxation

Like the US system, the Italian one is a marginal system with brackets, however in this case the fiscal pressure is higher and the brackets smaller. As a reference, it was used the 2020 Irpef brackets and percentages (Table 3.3). However, the brackets were calculated in dollars, so all the values were multiplied by 1.1575, the exchange rate at October 1st 2021.

Bracket in €	0-15000	15000-28000	28000-55000	55000-75000	75000+
Bracket in Coin	0-17	17-32	32-63	63-86	86+
Tax*	23%	27%	38%	41%	43%

Table 3.3: 2020 Irpef tax system: *this is the proportional system that was in place in the Italy in 2020, notice that the tax is marginal (i.e. if you are in the second bracket you will pay 3450€ + 27% of the amount above 15000€)*

In this case, there is an even higher pressure, which might push toward higher equality and lower total production.

3.3.4 Communism

Communism is the opposite of the Free market, in this simulation the fiscal pressure is 1, meaning that every 100 time-step all the agents are taxed for

the entire amount of coins they own and then redistributed. This system will probably deliver the lowest production and the highest equality.

3.4 Training Results

The process of training is computationally complex. It requires a good machine and plenty of time. In this case, there were used multiple machines and the trial period on Microsoft Azure allowed the use of up to 4 CPUs. No GPUs were implemented for the training purpose since it wasn't present in all the machines. However, by implementing an FCNet instead of a CNN the loss of efficiency the lack of GPU is not enormous.

The time needed to run a single training for 200M of time-steps was between 100/200h depending on the machine. This huge amount of time needed did not give any possibility to do hyper-parameter training of any sort to increase efficiency, nor to check the consistency of the results.

In Figure 3.4 we can see the full training process for these simulations. We could start to draw some conclusions from it, however, the values displayed are not representative of the true production and equality. We can see it in Figure 3.2 where the average production is around 1400, whereas in Figure A.4 the total actual production is above 2500. This is not just a right tale coincidence. This is an issue (in data collection) that I discovered once all the data was trained, thus it is impossible to recreate the correct graphs unless running all the training again (800+ hours).

To produce valid and comparable results, the weights obtained from the training were used to run 1000 simulations. With this data set, now is possible to generate reliable and comparable data.

The first comparisons that are presented are the histograms on Figure 3.5. In particular we have that the most productive economy is the Free market with an average production of 2751 ($\sigma = 161$), and at the same time is the one with the highest disparity, 0.437 ($\sigma = 0.014$). US economy is the second

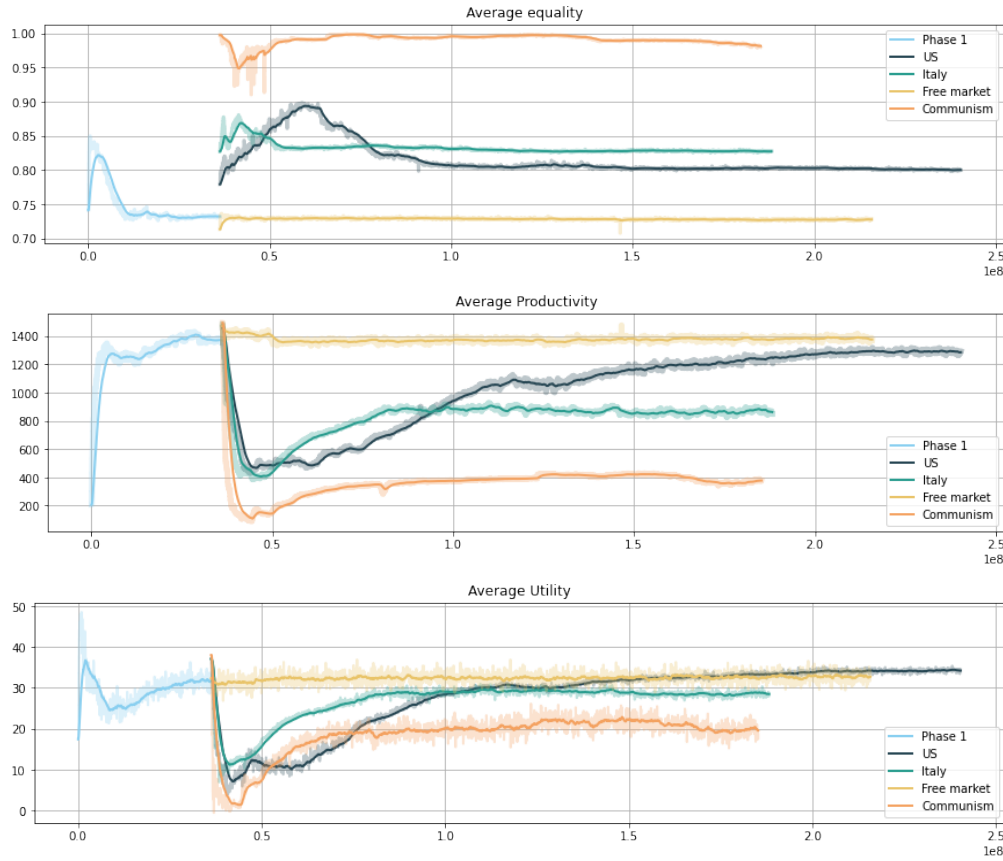


Figure 3.4: Phase 2 training brake down: *here we can see the training process, on the horizontal axis there are the time-steps, the training lasted for around 200M steps per simulation, and in the graphs is possible to see the evolution of equality, productivity and utility of the agents on average.*

most productive 2393 ($\sigma = 235$) and with equality 0.587 ($\sigma = 0.017$). Italy has a production of 1548 ($\sigma = 125$) and agents have a coin equality of 0.64 ($\sigma = 0.013$). And finally the Communism has the highest level of equality 0.99 ($\sigma = 0.009$) and the lowest total productivity 732 ($\sigma = 32$).

If we compare the results obtained so far with the one in the paper of Zheng and Trott, we can see that there is a loss in efficiency of -16.3% in Free market productivity and of 12% in equality. For the US instead, -10% and 22% respectively. This loss in overall efficiency was already expected as discussed above.

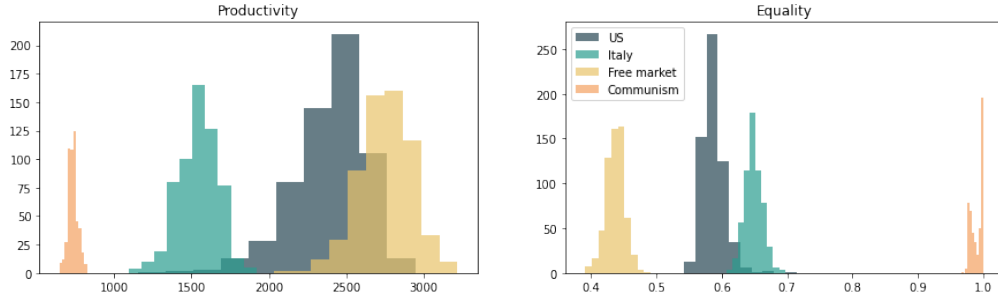


Figure 3.5: Productivity and Equality histogram: *these are the two histograms with the correct data on production (on the left) and equality (on the right). The color coding is the same across the two graphs.*



Figure 3.6: Coin Distribution: *this is a graphical representation of the coin distribution in the 4 experiments. The darker the color the richer the agent. The areas are proportional to the amount of economy owned by the agent.*

To have a better understanding of these results, we can check Figure 3.7. In this figure, equality is plotted in the plane against productivity. If we take into consideration the two systems that are actually used (US and Italy), we can see how the fiscal pressure given by Italian brackets and percentages causes a loss in productivity of about 35.4% compared to the US, but there is a gain of only 10.54% in equality.

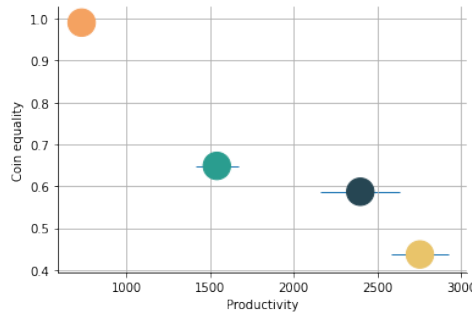


Figure 3.7: Coin Distribution

The last performance measure that could be taken into consideration is the product of the previous two. It makes sense, as a policymaker, to try to maximize not only the production or the equality but both of them at the same time as seen in economical theory [3]. Thus the usage of the product of the two measures seen so far could be a good candidate objective for the tax maker. In Figure 3.8 we can see the distribution of this value across the 4 experiments. Here is noticeable how the US economy shows a higher value compared to the Free market.

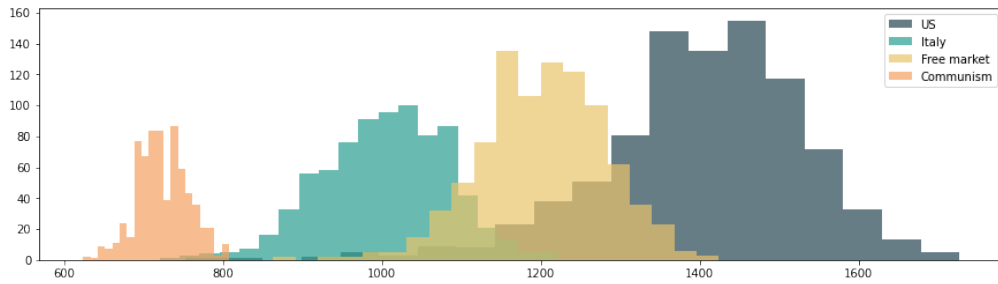


Figure 3.8: Coin equality times total production

These results, however, are not representative of the real impact of the actual taxation systems on the population. This simulation is still far away from reality. In addition, several necessary tests on consistency of the results and hyperparameter tuning are missing.

However, we can still conclude that AI-driven agents can perform tasks that otherwise would be impossible to model. The economical relevance of the results is now questionable, however, it could be interesting to further develop simulations where economical agents are subject to RL optimization. The use of such simulations could be a starting point to check whether the agents' behaviour converges to the analytical behaviour. Moreover, my work can be developed into more complex simulations to explore the effect of policies.

Conclusions

We have seen how reinforcement learning could be used to produce economical simulations. In detail, we saw RL applied to the gather-and-trade setup created by Alex Trott and Stephan Zheng. Here, we noticed how it is possible to maximize agents' utility through proximal policy optimization, a kind of gradient descent optimization. And we saw how this gave birth to agents' specialization, a feature that is hard to capture analytically. This first result is noticeable since it is a proof of concept that it is possible to create an experimental setup that is, once trained, easy to interact with and easy to run multiple times to gather data.

Afterward, we saw four different taxation systems (and redistribution) applied to a common starting point and commented on the difference in productivity and equality. When we compared the reproduction results, we noticed a loss in efficiency. This was justified from the model utilized and from the lack of hyperparameter tuning. Nonetheless, we reasonably found out, that there is a negative impact of taxation on production and a positive impact on equality. This impact cannot be properly quantified as a partial effect for a multitude of reasons. At first, there is the need for an efficient model, then, we need enough data to draw statistical conclusions.

Future work

This simulation can be a nice starting point for a multitude of future works. At first, it is possible to study the application of RL in a simpler setup that can

be modeled in an analytical way and check if RL converges to the analytical solutions. Afterward, starting from this model is possible to explore plenty of scenarios. For example, the brain drain and tax evasion could be interesting to model. For the former, we could simulate two countries that coexist and we could give the agents the choice to move from country A to country B. Then we could study how the high and low skilled agents decide to move across the countries. For the latter is possible to add the choice to reveal the real income and a small possibility of being caught. Here the interest would be in if or how the agents decide to evade given different tax systems.

These are just two examples, but the freedom of design is technically endless. Of course, such research would need better machines, thus computational power, than the one's used for this thesis.

Appendix A

Tables, Algorithms and Graphs

Parameter		Value
Episode Lenght	H	1000
World height		25
World width		25
Resources respawn prob.		0.01
Initial Coin endowment		0
Iso-elastic utility exponent	η	0.23
Move labor		0.21
Gather labor		0.21
Trade labor		0.05
Build labor		2.1
Base build payout		10
Max skill multiplier		3
Max bid/ask price		10
Max bid/ask order duration		50
Max simultaneous open orders		5
Tax period duration	M	100
Min bracket rate		0%
Max bracket rate		100%

Table A.1: Hyperparameters for the environment.

Parameter		Value
Training Algorithm		ppo
Number of Agents	N	4
Number of parallel environment replicas		60
Sampling horizon (steps per replica)	h	200
Train batch size		3000
SGD minibatch size		500
Number SGD iter		10
CPUs		2
Learning rate	Lr	3e-4
Learning rate schedule (phase 1)		[[0, Lr], [10M, Lr], [30M, 0]]
Learning rate schedule (phase 2)		[[35M, Lr], [50M, 1e-6]]
Entropy regularization coefficient		1e-4
Gamma	γ	0.99
GAE lambda		0.9
Gradient clipping parameter		0.25
Value function loss coefficient		0.05
Number of fully-connected layers		2
Agents get full spatial observation		False
Agent spatial observation box half-width		5
Phase one training duration		35M
Phase one energy warm-up constant	ϑ	10000
Phase two training duration		150M

Table A.2: Hyperparameters for the Training.

Variable Name	Dimension
world Map	(7, 11, 11)
world-idx_map	(2, 11, 11)
world-loc-row	(1,)
world-loc-col	(1,)
world-inventory-Coin	(1,)
world-inventory-Stone	(1,)
world-inventory-Wood	(1,)
time	(1, 1)
Build-build_payment	(1,)
Build-build_skill	(1,)
ContinuousDoubleAuction-market_rate-Stone	(1,)
ContinuousDoubleAuction-price_history-Stone	(11,)
ContinuousDoubleAuction-available_asks-Stone	(11,)
ContinuousDoubleAuction-available_bids-Stone	(11,)
ContinuousDoubleAuction-my_asks-Stone	(11,)
ContinuousDoubleAuction-my_bids-Stone	(11,)
ContinuousDoubleAuction-market_rate-Wood	(1,)
ContinuousDoubleAuction-price_history-Wood	(11,)
ContinuousDoubleAuction-available_asks-Wood	(11,)
ContinuousDoubleAuction-available_bids-Wood	(11,)
ContinuousDoubleAuction-my_asks-Wood	(11,)
ContinuousDoubleAuction-my_bids-Wood	(11,)
Gather-bonus_gather_prob	(1,)
action_mask	(50,)

Table A.3: Full observation space.

Algorithm 2 Agents Learning Loop.

Require: Sampling horizon h and tax period M **Require:** On-policy learning algorithm \mathcal{A} ▷ PPO**Require:** Stopping criterion ▷ i.e. utility convergence**Ensure:** Trained agent policy weights θ $s, \mathbf{o} \leftarrow s_0, \mathbf{o}_0$ ▷ reset episode $\theta \leftarrow \theta_0$ $D \leftarrow \{\}$ ▷ reset transition buffer**while** training **do** **for** $t = 0, 1, \dots, h$ **do** $\mathbf{a} \leftarrow \pi(\cdot | \mathbf{o}, \theta)$ $s', \mathbf{o}', \mathbf{r} \leftarrow \text{Env.step}(s, \mathbf{a}, \tau)$ **if** $t \bmod M = M - 1$ **then** $s', \mathbf{o}', \mathbf{r} \leftarrow \text{Env.tax}(s')$ **end if** $D \leftarrow D \cup \{(\mathbf{o}, \mathbf{a}, \mathbf{r}, \mathbf{o}')\}$ $s, \mathbf{o} \leftarrow s', \mathbf{o}'$ **end for** Update θ using data in D and \mathcal{A} $D \leftarrow \{\}$ **if** episode is completed **then** $s, \mathbf{o} \leftarrow s_0, \mathbf{o}_0$ ▷ reset episode **end if** **if** the stopping criterion has been met **then return** θ **end if****end while**

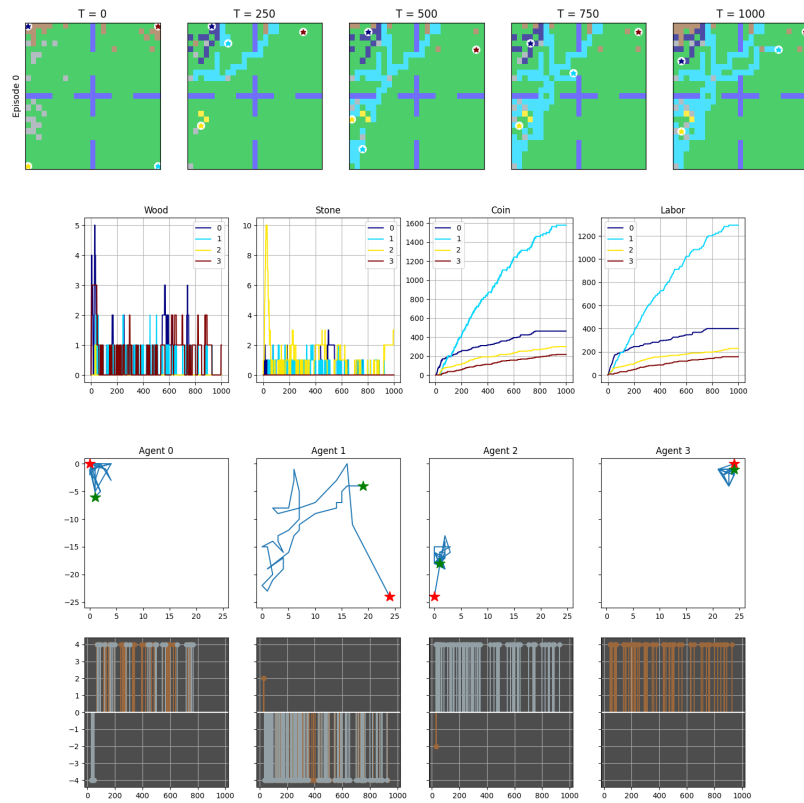


Figure A.1: Free market full brake down: *this is a more complete overview of one random simulation crated using the weights at the end of phase 2 training for the Free market.*

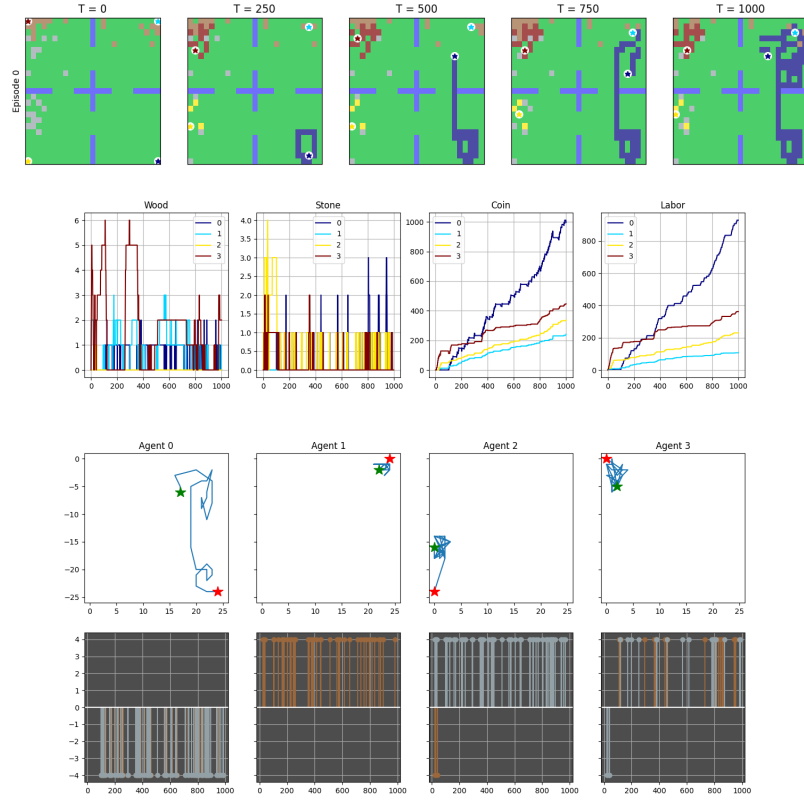


Figure A.2: US full brake down: *this is a more complete overview of one random simulation crated using the weights at the end of phase 2 training for US taxation.*

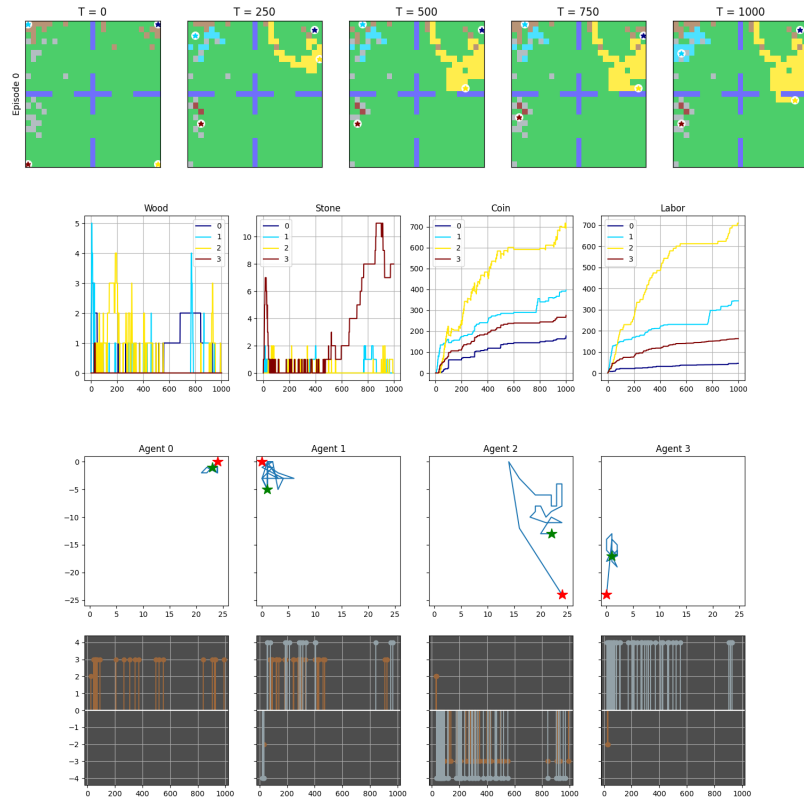


Figure A.3: Italy full brake down: *this is a more complete overview of one random simulation crated using the weights at the end of phase 2 training for Italian taxation.*

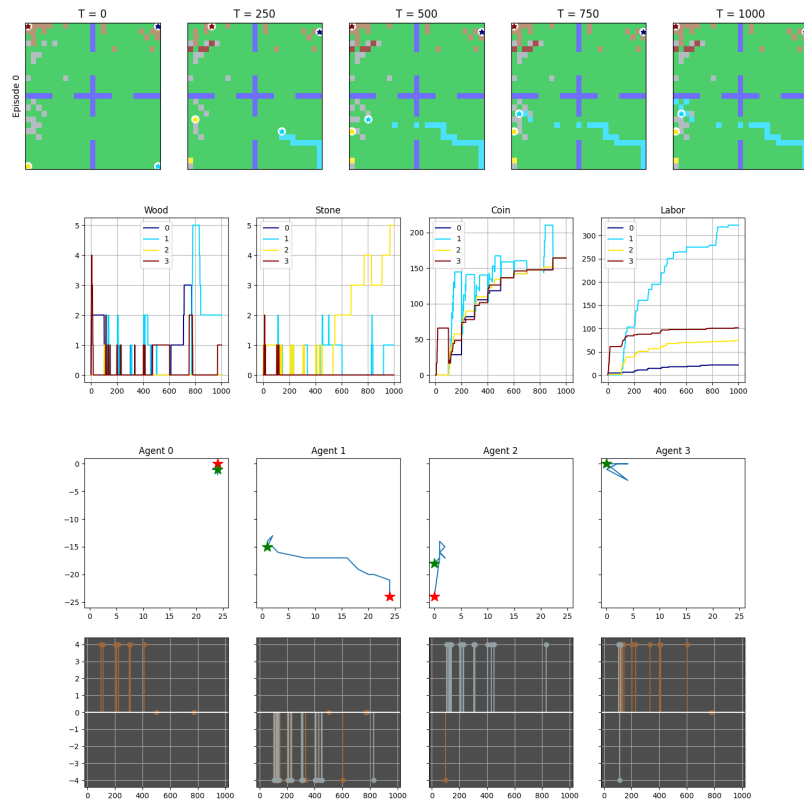


Figure A.4: Communism full brake down: *this is a more complete overview of one random simulation crated using the weights at the end of phase 2 training for communism.*

Bibliography

- [2] Stephan Zheng et al. “The AI Economist: Optimal Economic Policy Design via Two-level Deep Reinforcement Learning”. In: *arXiv preprint arXiv:2108.02755* (2021).
- [3] Emmanuel Saez. “Using elasticities to derive optimal income tax rates”. In: *The review of economic studies* 68.1 (2001), pp. 205–229.
- [4] Stephan Zheng et al. “The ai economist: Improving equality and productivity with ai-driven tax policies”. In: *arXiv preprint arXiv:2004.13332* (2020).
- [5] Richard S Sutton and Andrew G Barto. “Reinforcement learning: An introduction”. In: (2018).
- [6] John Schulman et al. “High-dimensional continuous control using generalized advantage estimation”. In: *arXiv preprint arXiv:1506.02438* (2015).
- [7] John Schulman et al. “Proximal policy optimization algorithms”. In: *arXiv preprint arXiv:1707.06347* (2017).
- [8] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [9] Anthony B Atkinson and Joseph E Stiglitz. “Lectures on public economics: Updated edition”. In: (2015).

Sitography

- [1] *ai-economist PyPI*. URL: <https://pypi.org/project/ai-economist>.

- [10] *The AI Economist - Salesforce blogpost*. URL: <https://einstein.ai/the-ai-economist>.
- [11] *AI Economist - Github*. URL: <https://github.com/salesforce/ai-economist>.
- [12] *Akira Sosa - Github*. URL: <https://github.com/akirasosa>.
- [13] *Akira Sosa - Medium*. URL: <https://medium.com/vitalify-asia/ai-economist-sums-up-communism-21260e5540f8>.

Source Code

- [14] *Lorenzo Mezzini - GitHub*. URL: https://github.com/lorenzomezzini/master_thesis.