

The AI economist: Improving Equality and Productivity with AI-Driven Tax Policies

Stephan Zheng, Alexander Trott, Sunil Srinivasa, Nikhil Naik, Melvin Gruesbeck, David C. Parkes, and Richard Socher, 2020, mimeo.

Presenter: Yoji Tomita

RL-GT ㊦ June 17, 2021

Table of Contents

1. Introduction
2. Economic Simulations: Learning in Gather-and-Build Games
 - 2.1 Notation and Preliminaries
 - 2.2 Environment Rules and Dynamics
 - 2.3 Using Machine Learning to Optimize Agent Behavior
3. Machine Learning for Optimal Tax Policies
 - 3.1 Periodic Taxes with Bracketed Schedules
 - 3.3 Inner-Outer-Loop Reinforcement Learning
4. Improved Social Outcomes with AI Agents
 - 4.2 Training Strategy: Two-phase Training and Tax Curricula
 - 4.3 Equality, Productivity, and Social Welfare Metrics
 - 4.4 Tax Schedules and Wealth Redistribution after Taxes and Subsidies

1. Introduction

- ・ イントロダクション

2. Economic Simulations: Learning in Gather-and-Build Games

- Economic environment について.
- まずは税の無い設定 ("free-market") で説明する.

2.1 Notation and Preliminaries

- Partial-observable multi-agent Markov Games(MGs): $(S, A, r, \mathcal{T}, \gamma, o, \mathcal{I})$
 - S : 状態空間 (state space)
 - A : 行動空間 (action space)
 - $r_{i,t}$: 報酬関数 (reward function)
 - \mathcal{T} : 遷移関数 (transition function) $s_{t+1} \sim \mathcal{T}(\cdot \mid s_t, \mathbf{a}_t)$
 - γ : 割引因子 (discount factor)
 - $o_{i,t}$: 観測 (observation)
 - time step $t = 0, 1, \dots, H$.

- Agents' policy : $\pi_i(\cdot \mid o_{i,t}, h_{i,t}; \theta_i)$
 - $h_{i,t}$: hidden state (自分の私的情報と, 過去の history)
 - θ_i : policy の parameter
 - エージェント i は次の最大化問題を得く policy を求める:

$$\max_{\theta_i} \mathbb{E}_{a_i \sim \pi_i, \mathbf{a}_{-i} \sim \boldsymbol{\pi}_{-i}, s' \sim \mathcal{T}} \left[\sum_t \gamma^t r_{i,t} \right]. \quad (1)$$

- データ効率性のため, すべてのエージェントは training の間パラメータ θ を共有する.
- 彼らの行動 $\pi_i(a_i \mid o_i, h_i; \theta)$ は, agent-specific observations o_i と hidden-state h_i によって異なる.

t	time
i, j, k	agent indices
θ, ϕ	model weights
s	state
o	observation
a	action
r	reward
π	policy
γ	discount factor
\mathcal{T}	state-transition, world dynamics
h	hidden state

x	endowment
x^c	coin
x^s	stone
x^w	wood
z	income
l	labor
u	utility
T	tax
τ	tax-rate
π_p	planner policy
swf	social welfare
ω	social welfare weight
g	social marginal welfare weight
gini	Gini index
eq	Equality index

Table 1: Notation. Subscripts are indices. Superscripts are labels.

2.2 Environment Rules and Dynamics

Gather-and-Build game

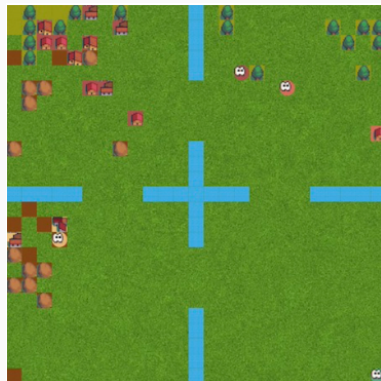
- ・ 2次元の grid (25×25) からなる世界が舞台.
- ・ エージェントはフィールドを歩き回り, 資源 (石と木) を集め, それらを1つずつ使って家を建て, また資源を coin を介してトレードする.
- ・ 資源は空タイルに確率的に産み出される.
- ・ エージェントは家を建てると coin が得られるが, 得られる coin は agent の skill ごとに異なる.

Labor and Skill.

- Agent の action(moving, gathering, trading, building) にはそれぞれ labor cost が設定されている.
- 各 time に agent がどれか 1 つ行動をとると, その labor cost がかかる.
- building skill (1 以上 3 以下) が各 agent に設定されていて, 家を建てると agent は $10 \times \text{skill}$ 分の coin を得る.
- collection skill (1 以上 2 以下) もあり, 資源を拾うとこの skill 分の資源を得る (skill 1.2 の場合, 確定で 1 つ資源を得て, さらに確率 0.2 でもう 1 つ資源を得る)

Environment Scenario.

- field は水により 4 つの区域に別れている（水部分は通れない）
- 資源は空間的に集まって発生する.
- 4 agents
- building skills は 1.13, 1.33, 1.65, 22.2（Pareto 分布 w/ exponent $a = 4$, scale $m = 1$ の quartiles を元に設定）
- 1 episode は $H = 1000$ time steps からなる.



2.3 Using Machine Learning to Optimize Agent Behavior

- Agent の utility function:

$$u_i(x_{i,t}, l_{i,t}) = \text{ccra}(x_{i,t}^c) - l_{i,t}, \quad \text{where } \text{ccra}(z) = \frac{z^{1-\eta} - 1}{1-\eta}, \quad \eta > 0. \quad (2)$$

- $x_{i,t} = (x_{i,t}^w, x_{i,t}^s, x_{i,t}^c)$: i の保有する木・石・コイン.
- $l_{i,t}$: 蓄積労働量.
- η : エージェントの utility function の non-linearity をコントロールするパラメータ.
- Rational economic agent は以下の最大化を行う.

$$\forall i : \max_{\pi_i} \mathbb{E}_{a_i \sim \pi_i, \mathbf{a}_{-i} \sim \pi_{-i}, s' \sim \mathcal{I}} \left[u_i(x_{i,0}, l_{i,0}) + \sum_{t=1}^H \gamma^t \underbrace{(u_i(x_{i,t}, l_{i,t}) - u_i(x_{i,t-1}, l_{i,t-1}))}_{=r_{i,t}} \right]. \quad (3)$$

Deep RL agents

- deep neural network を用いる agent policy を modelling する:

$$a_{i,t} \sim \pi(o_{i,t}^{\text{world}}, o_{i,t}^{\text{agent}}, o_{i,t}^{\text{market}}, o_{i,t}^{\text{tax}}, h_{i,t-1}; \theta)$$

- $o_{i,t}^{\text{world}}$: 近くの状況に関する観測.
- $o_{i,t}^{\text{agent}}$: public な agent の状況 (資源・コイン保有) と, private agent states(skill 値と labor performed)
- $o_{i,t}^{\text{market}}$: transfer market の状況 (bid, ask offer)
- $o_{i,t}^{\text{tax}}$: tax rates

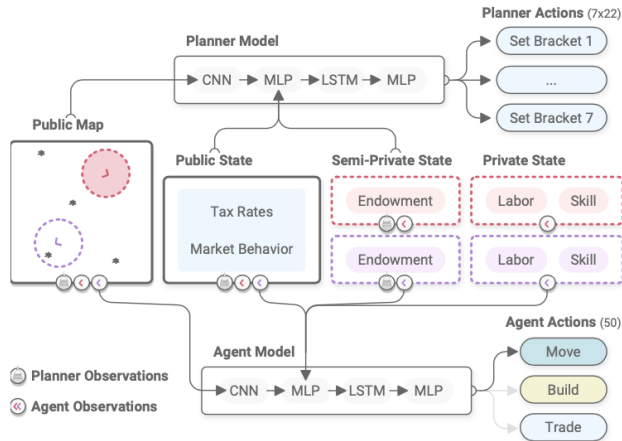
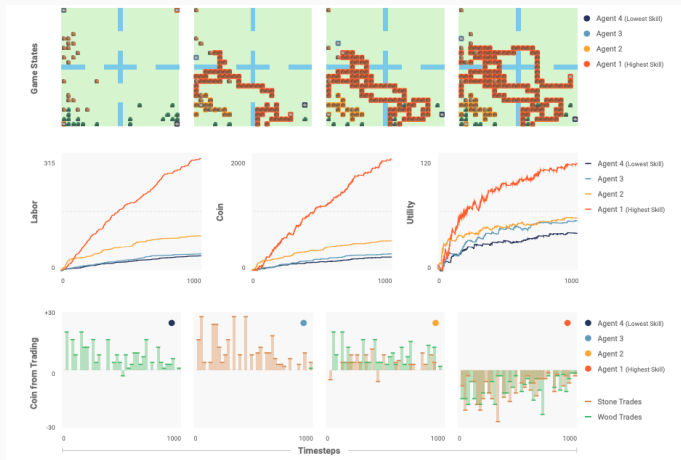


Figure 3: Schematic overview of the general network architecture used in our work. Spatial observations are processed by a stack of two convolutional layers (CNN) and flattened into a fixed-length feature vector. This feature vector is concatenated with the remaining observation inputs and the result is processed by a stack of two fully connected layers (MLP). The output is then used to update the hidden state of an LSTM and action logits are computed via a linear projection of the updated hidden state. Finally, the network computes a softmax probability layer for each action head. For the agent policy, there is a single action space and action head. For the tax policy, there is a separate action space and action head for each tax rate the tax policy controls (described below).

Emergent Behavior of AI Agents



- 左図は no-tax 下で train 後の AI agents の 1 episode の行動の例.
- low-skill agents (紺, 水) は資源を集めて market で売ることには徹している.
- high-skill agent (オレンジ) は market で資源を買って家を立てている.
- 黄色は最初は家を立ててるが, のちに資源を得る方にスイッチしている.

3. Machine Learning for Optimal Tax Policies

- ・ 課税と再分配を行う social planner を導入する.
- ・ Social planner は, 生産性と平等性の trade-off に直面している.
 - ・ 無課税 (free-market) では生産性は最大化されるが, 不平等.
 - ・ 課税・再分配を行うと平等性が増すが, 生産性が落ちる.
- ・ ここでは, free-market, US-federal, Saez framework, AI economist による social planner を試す.

3.1 Periodic Taxes with Bracketed Schedules

Income Taxes.

- Tax period は M steps 続く（実験では $M = H/10$ とし, 1 episode に 10 tax periods があるものとする）
- ピリオド p の税は, time step t から $t + M$ までの収入 z_i^p に課される.
- Tax period の初めに, social planner は tax schedule $T(z)$ を決めて公表する.
 - 各 agent i は, 収入 z_i^p に応じて $T(z)$ を支払う.
 - 集められた税は, 全 agent に平等に分配される.
 - よって, 分配後の agent i の収入は,

$$\tilde{z}_i^p = z_i^p - T(z_i^p) + \frac{1}{N} \sum_{j=1}^N T(z_j^p). \quad (5)$$

Bracketed Tax Schedules.

- Scheme 間の比較を可能にするため, tax schedule は次のように "bracketed" されたもののみを考える.
- Cut-off income levels $\{m_b\}_{b=0}^B$ s.t. $0 = m_0 \leq m_1 \leq \dots \leq m_{B-1} \leq m_B = +\infty$ が先に与えられている.
- Social planner は, 各 bracket b に含まれる収入に対して適用される marginal tax rate $\tau_b \in [0, 1]$ を選ぶことで, tax schedule $T(\cdot)$ を決定する.

$$T(z) = \sum_{b=0}^{B-1} \tau_b \cdot ((m_{b+1} - m_b) \cdot 1[z > m_{b+1}] + (z - m_b) \cdot 1[m_b < z \leq m_{b+1}]) .$$

3.2 Optimal Taxation

Social Welfare Functions

- Social planner の目的関数である social welfare function は, 生産性と平等性の trade-off を組み込めるように次のように決める.
- エージェントのコイン保有 $x^c = (x_1^c, \dots, x_N^c)$ に対し, equality を次で定義:

$$\mathbf{eq}(x^c) = 1 - \mathbf{gini}(x^c) \cdot \frac{N}{N-1}, \quad 0 \leq \mathbf{eq}(x^c) \leq 1. \quad (7)$$

where

$$\mathbf{gini}(x^c) = \frac{\sum_{i=1}^N \sum_{j=1}^N |x_i^c - x_j^c|}{2N \sum_{i=1}^N x_i^c}, \quad 0 \leq \mathbf{gini}(x^c) \leq \frac{N-1}{N} \quad (8)$$

- \mathbf{eq} は, 1 で完全に平等 (全員同じ収入), 0 で完全に不平等 (1 人が全コインを独占).

- ・ 生産性は,

$$\mathbf{prod}(\mathbf{x}^c) = \sum_{i=1}^N x_i^c. \quad (9)$$

- ・ この **eq** と **prod** を social welfare function とする.¹

$$\mathbf{swf}_t(\mathbf{x}_t^c) = \mathbf{eq}_t(\mathbf{x}_t^c) \cdot \mathbf{prod}_t(\mathbf{x}_t^c). \quad (10)$$

¹Social welfare function として, weight $\omega_i \geq 0$ を用いて

$$\mathbf{swf}_t(\mathbf{x}_t^c, \mathbf{l}_t) = \sum_{i=1}^N \omega_i \cdot u_i(x_{i,t}^c, l_{i,t}). \quad (11)$$

を用いることも可能.

The Planner's Problem.

- Social Planner は,
 - agent の保有資源・コイン $x_{i,t}$, フィールドの状態 (agent・資源の位置) と market の状況は観測可能,
 - 各 agent の skill は直接には観測できない.
- Planner の最大化問題は,

$$\max_{\pi_p} \mathbb{E}_{\tau \sim \pi_p, \mathbf{a} \sim \pi, s' \sim \mathcal{T}} \left[\mathbf{swf}_0 + \sum_{t=1}^H \gamma^t \underbrace{(\mathbf{swf}_t - \mathbf{swf}_{t-1})}_{=r_{p,t}} \right]. \quad (12)$$

3.3 Inner-Outer-Loop Reinforcement Learning

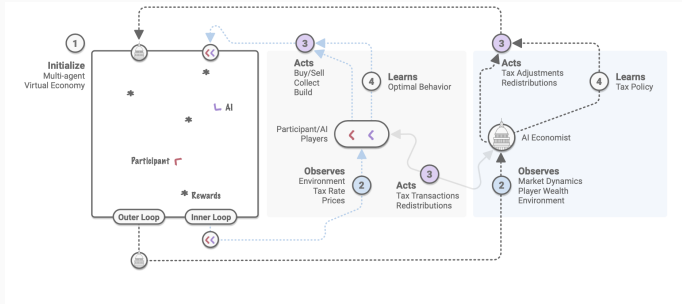


Figure 6: Two-level RL. In the inner loop, RL agents gain experience by performing labor, receiving income, and paying taxes, and learn through balancing exploration and exploitation how to adapt their behavior to maximize their utility. In the outer loop, the social planner adapts tax policies to optimize its social objective.

- ・この状況では, Inner loop で agents が行動を選んで学習し, Outer loop で planner が税制を選んで学習することになる.
- ・Inner loop での Agent の学習行動と, outer loop で planner の選ぶ税制が相互に影響し合うため, 報酬が agent と planner 双方にとって不安定になる.

Algorithm 1 Inner-Outer Loop Reinforcement Learning. Economic agents and social planner learn simultaneously. Bold-faced symbols indicate quantities for multiple agents. Note that agents share weights.

Require: Sampling horizon \bar{h} , tax period length M

Require: On-policy learning algorithm \mathbb{A} (for instance, A3C, PPO)

Require: Stopping criterion C (for instance, agent and planner rewards have not improved)

Ensure: Trained agent and planner policy weights θ, ϕ

$s, \mathbf{o}, \mathbf{o}_p, \mathbf{h}, h_p \leftarrow s_0, \mathbf{o}_0, \mathbf{o}_{p,0}, \mathbf{h}_0, h_{p,0}$ ▷ Reset episode

$\theta, \phi \leftarrow \theta_0, \phi_0$ ▷ Initial agent and planner policy weights

$D, D_p \leftarrow \{\}, \{\}$ ▷ Reset agent and planner transition buffers

while training **do**

for $t = 1, \dots, \bar{h}$ **do**

$\mathbf{a}, \mathbf{h} \leftarrow \pi(\cdot | \mathbf{o}, \mathbf{h}, \theta)$ ▷ Sample agent actions; update hidden state

if $t \bmod M = 0$ **then** ▷ First timestep of tax period

$\tau, h_p \leftarrow \pi_p(\cdot | \mathbf{o}_p, h_p, \phi)$ ▷ Sample marginal tax rates; update planner hidden state

else

 no-op, $h_p \leftarrow \pi_p(\cdot | \mathbf{o}_p, h_p, \phi)$ ▷ Only update planner hidden state

end if

$s', \mathbf{o}', \mathbf{o}'_p, \mathbf{r}, r_p \leftarrow \text{Env.step}(s, \mathbf{a}, \tau)$ ▷ Next state / observations, pre-tax reward, planner reward

if $t \bmod M = M-1$ **then** ▷ Last timestep of tax period

$s', \mathbf{o}', \mathbf{o}'_p, \mathbf{r}, r_p \leftarrow \text{Env.tax}(s', \tau)$ ▷ Apply taxes; compute post-tax rewards

end if

$D \leftarrow D \cup \{(\mathbf{o}, \mathbf{a}, \mathbf{r}, \mathbf{o}')\}$ ▷ Update agent transition buffer

$D_p \leftarrow D_p \cup \{(\mathbf{o}_p, \tau, r_p, \mathbf{o}'_p)\}$ ▷ Update planner transition buffer

$s, \mathbf{o}, \mathbf{o}_p \leftarrow s', \mathbf{o}', \mathbf{o}'_p$

end for

 Update θ, ϕ using data in D, D_p and \mathbb{A} .

$D, D_p \leftarrow \{\}, \{\}$ ▷ Reset agent and planner transition buffers

if episode is completed **then**

$s, \mathbf{o}, \mathbf{o}_p, \mathbf{h}, h_p \leftarrow s_0, \mathbf{o}_0, \mathbf{o}_{p,0}, \mathbf{h}_0, h_{p,0}$ ▷ Reset episode

end if

if criterion C is met **then return** θ, ϕ

end if

end while

4. Improved Social Outcomes with AI Agents

4.1 Baseline Methods

- 次の4つの tax model を比較する.
 - free-market (no taxes)
 - US federal single-filer 2018 tax schedule
 - Saez tax formula (adapted for a multi-period setting)
 - AI Economist planner
- Tax bracket は, 全ての tax model に共通で, 2018 US federal incom tax をもとに 1/1000 にスケーリングした以下を用いる.

$$\mathbf{m} = [0, 9.7, 39.475, 84.2, 160.725, 204.100, 510.3, \infty]. \quad (15)$$

US Federal Income Tax Rates (Single-filer, 2018).

- 2018 US federal income tax をもとに, bracket tax rates は,

$$\tau = [0.1, 0.12, 0.22, 0.24, 0.32, 0.35, 0.37] \quad (16)$$

とする.

Saez Tax Formula (single-period)

- Saez [2001] をもとに, まず single-period economy での optimal tax rates を求める.
- f, F は (pre-tax) 収入の分布の probability density と cumulative distribution function とし, planner はそれらを観測できるとする.
- Saez [2001] では, まず linear-weighted social welfare functions (11) に対し, social marginal welfare weights

$$g_i = \frac{d\mathbf{swf}}{du_i} \frac{du_i}{dx_i^c} = \omega_i \frac{du_i}{dx_i^c}.$$

を求める.

- これを我々のモデルに当てはめるために $g_i = \frac{1}{z_i}$ とし, さらにこれを normalize してよい $\sum_{i \in \mathcal{I}} g_i = 1$.

- $\alpha(z)$ を, the marginal average income at income z , normalized by the fraction of incomes above z , つまり,

$$\alpha(z) := \frac{z \cdot f(z)}{1 - F(z)}. \quad (18)$$

- $G(z)$ は, normalized, reverse cumulative Pareto weight over incomes above a threshold z

$$G(z) := \frac{1}{1 - F(z)} \int_{z'=z}^{\infty} f(z')g(z')dz'. \quad (19)$$

- また, 弾力性 (elasticity) $e(z)$ を, average sensitivity of an agent's income to changes in the tax rate, つまり

$$e(z) = \frac{dz/z}{d(1 - \tau(z))/(1 - \tau(z))}. \quad (20)$$

- Saez [2001] によると, optimal marginal tax-rate は

$$\tau(z) = \frac{1 - G(z)}{1 - G(z) + \alpha(z)e(z)} \quad (21)$$

となり, これは income distribution と elasticity に依存して決まる.

Saez Tax Formula (multi-period)

- Saez formula (21) を使うには, income elasticity $e(z)$ の推定が必要.
- Gruber and Saez [2002] に従い, constant tax elasticity \tilde{e} を仮定すると,

$$z_t = z^0 \cdot (1 - \tau_t)^{\tilde{e}}. \quad (22)$$

- したがって,

$$\log(z_t) = \tilde{e} \cdot \log(1 - \tau_t) + \log(z^0) \quad (23)$$

を過去の tax period のデータから OLS で推定して用いる.

- AI Economist は, deep neural network により各 bucket の marginal tax rate を決める.

$$\tau \sim \pi_p \left(o_{p,t}^{\text{world}}, o_{p,t}^{\text{agent}}, o_{p,t}^{\text{market}}, o_{p,t}^{\text{tax}}, h_{p,t-1}; \phi \right). \quad (24)$$

- Basic network architecture は agent と同様 (Figure 3).
- agent と planner は持つ observation が異なる (planner はフィールド全体が見れるが, agent の private skill は見れない) .

4.2 Training Strategy: Two-phase Training and Tax Curricula

- Section 3.3 で述べたように, この joint optimization は不安定性をもたらす.
- 学習を安定化させるため, ここでは次の two-phase training approach を行った.
 - First phase: agent models の集合に対し, 無課税 (no-tax scenario) で学習させる.
 - Second phase: 税モデルを融合にし, エージェント・プランナーの学習を継続する. また突然の税導入による不安定性を回避するため, 限界税率の上限を 10% から 100% に徐々に引き上げる.
- また, planner policy に entropy regularization を行うことも効果的だった.
 - Entropy regularization は, policy gradient objective に次の additional, weighted term

$$\text{entropy}(\pi) = -\mathbb{E}_{a \sim \pi(\cdot | s)} [\log \pi(a | s)].$$

を加える (Williams and Peng [1991], Mnih et al. [2016]) .

- 実験は RLlib framework [Liang et al., 2018] を用いて実施.
- また Policy gradients の計算に, proximal policy gradients [Schulman et al., 2017] と Adam optimizer [Kingma and Ba, 2014] を使う.
- サンプルは 60 の environment replicas から平行して集め, sampling horizon は 200 timesteps とする.
- 全ての実験で, 400 million sample の phase two training を実施し, これは agent と planner model は stable policy への収束に十分であった.

4.3 Equality, Productivity, and Social Welfare Metrics

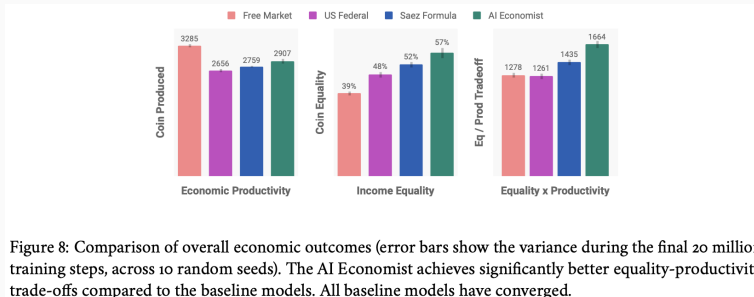


Figure 8: Comparison of overall economic outcomes (error bars show the variance during the final 20 million training steps, across 10 random seeds). The AI Economist achieves significantly better equality-productivity trade-offs compared to the baseline models. All baseline models have converged.

- Figure 8 は episode 終了時の各 tax model での economic outcomes の比較.
- Tax は常に productivity を下げるが, その下り幅は AI Economist が最も少ない.
- Income equality ($1 - \text{Gini index}$) は, AI economist で最も高い.
- Equality と Productivity の積も AI economist が最も高く, 次に高い Saez formula と比べても 16% 高い.

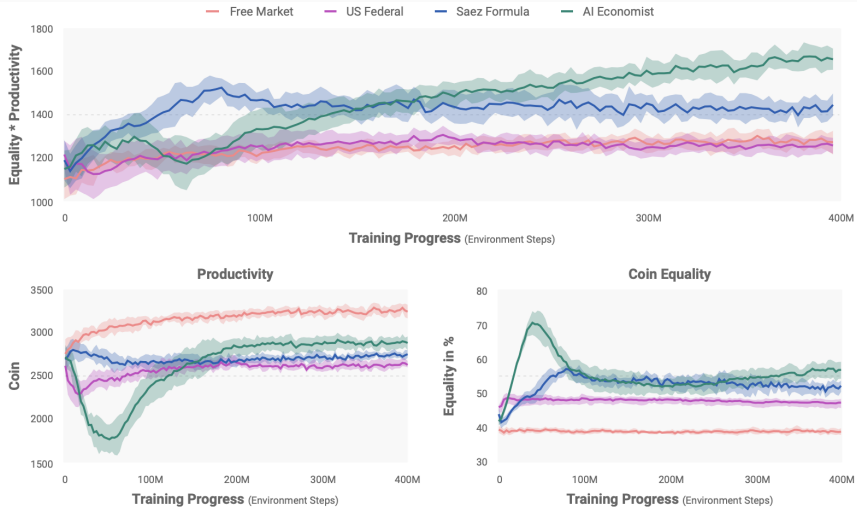


Figure 7: Empirical training progress for all models. The AI Economist (Green) achieves significantly better social outcomes than the baseline models. All baseline models have converged.

4.4 Tax Schedules and Wealth Redistribution after Taxes and Subsidies

Comparing Tax Schedule

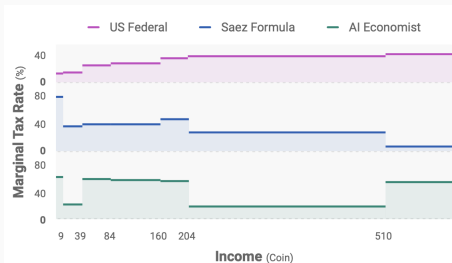


Figure 9: Comparison of average tax rates per episode. Variances within the Saez and AI Economist schedules are not shown. On average, the AI Economist sets a higher top tax rate than both of the US Federal and Saez tax schedules. The free-market collects zero taxes.

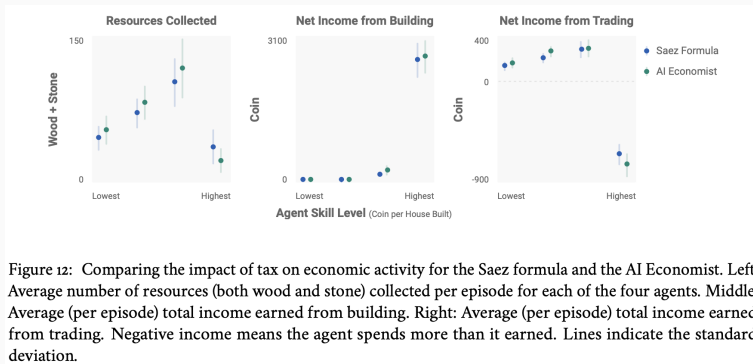
- Figure 9 は, 各 tax model での marginal tax rate の比較.
- US federal は, higher income に対して marginal tax rate は上がっていく.
- Saez formula では概ね逆.
- AI Economist はその blend になっている.

Effective Tax After Redistribution.



Figure 11: Agent-by-agent averages after sorting by skill. Income before redistribution (top-left) shows the average pre-tax income earned by each kind of agent. The amount of tax paid before distribution is shown in the bottom left. The amount of tax paid after redistribution is shown in the bottom right (the lower skill agents receive a net subsidy). The income after redistribution (top-right) shows the net average coin per agent at the

The Impact of Tax on Economic Activity



- Figure 12 は, Saez framework と AI economist での agent の行動の比較.
- AI economist では, high-skill worker はより資源を集めることより買って building することに特化し, low-skill workers はより資源を集めるようになっている.

4.5 Tax-Gaming Strategies

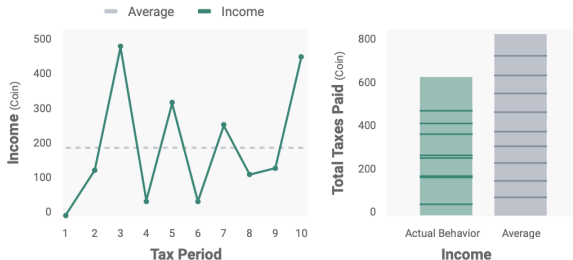


Figure 13: Left: Income of the highest-skilled agent for each tax period in an example episode with the AI Economist (green line). The dashed grey line shows the agent's average income. Right: Comparison of the total amount of tax the agent owed based on its actual income (green) and the tax it would have owed if it reported its average income in each period (gray). Each box in the column denotes the tax obligation in a single period.

- Figure 13 は AI economist での high-skilled worker の 1 episode の income の推移。
- 高収入の tax period と低収入の tax period が交互に起きている。
- Agent は収入をばらつかせることで課税額を抑えるように行動している (Saez でも)