

Evolution of Labour Market Mismatch through Adaptive Learning with Experience

Siting Estee Lu*

May 20, 2025

This is a PRELIMINARY draft.
[\[CLICK HERE FOR LATEST VERSION\]](#)

Abstract

Workers' past application choices can act as heuristics for future decisions. By integrating learning theory into search model, this work explores the role of experiences on workers' application choices. It provides an evolutionary perspective to the labour market dynamics, and offers some insights on equilibrium selection. I propose two market structures, where wages are unobservable and observable before workers make application decisions, and I model workers' learning behaviour using reinforcement learning and best response dynamics respectively. I show that in presence of multiple equilibria, experience-based learning generally leads to more efficient and asymptotically stable outcome of workers coordinating on applying with high probability to different firms in both static and dynamic wage settings. However, learning is nuanced as wages vary. Workers' choice probabilities are path-dependent and outcome may be more random. There can also be large fluctuations. Experience could serve as both an accelerator and inhibitor for reaching efficient outcome. Learning models not only highlight potential mechanism towards equilibrium, they also provide novel avenue for policies to improve market efficiency by exploring experience effect and intervention on the learning trajectory.

JEL: C73, D83, J64

Keywords: Job Search, Coordination, Adaptive Learning

*School of Economics, The University of Edinburgh. Contact: s.lu@ed.ac.uk.
I am indebted to Ed Hopkins and Axel Gottfries for their invaluable suggestions and continuous guidance. I would also like to express my gratitude to Michael Woodford for his valuable feedback and suggestions. The idea was previously presented at Toulouse Summer School in Quantitative Social Sciences 2023 and SGPE PhD Conference 2023.

1 Introduction

Workers' job application decisions could be guided by their past experiences, and modelling their learning behaviour offers new perspective on labour market dynamics. While experience-based adaptive learning have been studied in extensive-form games (Roth and Erev (1995)), normal-form games (Camerer and Hua Ho (1999)), choice predictions (Erev et al. (2010)), policy attitudes and future decisions (Albarracin and Wyer Jr (2000)), its application to search behaviour has not been thoroughly explored. Instead of imposing restrictions on workers' search strategies (Wright et al. (2021)), learning models provide viable mechanisms to substantiate equilibrium selection in presence of multiple equilibria, which could have important policy implications. They also suggest possible cognitive mechanisms underlying workers' search behaviour and shed light on several empirical puzzles, such as workers' apparent resistance to apply to higher wage jobs despite their ability to transit (Archer (2016)), and sorting at application stage due to cultural stereotypes associated with jobs (Barbulescu and Bidwell (2013)). Therefore, this paper investigates the role of past experiences on workers' adaptive learning behaviour in job applications and explores how this affects labour market outcome.

In this paper, I explore a framework of 2 workers and 2 firms. I propose two market structures, where wages are unobservable and observable to workers ex-ante to their application decisions respectively. These could offer some insights about the impact of experiences on learning dynamics and equilibrium behaviour under different market designs.

In the first market structure, workers learn solely based on feedback of jobs they have applied to. This process involves them recalling information from past "self" or learning through knowledge diffusion from prior generations, and revising their strategies via reinforcement learning (Erev and Roth (1998)). When firms set fixed wages, I found that workers learn to apply to different firms in the long run, and such pure strategy equilibria are asymptotically stable. Initial experiences could have a persistent impact on later choices and determine which equilibrium is more likely to be selected. As experiences may also create inertia to achieve more efficient one-to-one matching outcome, leading to longer learning trajectory, policymakers could consider reducing initial bias and lowering the impact from past memories to mitigate locked-in effect.

I also explore two-sided learning, where firms are allowed to adjust wages, similarly through reinforcement learning. In the long run, firms would drive wages down to 0 to maximize profit as workers' strategies stabilize, leading to a continuum of equilibria. Workers' equilibrium choice probabilities are path-dependent and can be more random. Along the learning trajectory, as wages are dynamically changing, workers could be navigating through learning about different sets of Nash equilibria over time. It could be harder to mitigate mismatch due to substantial accumulation of past experiences in payoff structure with a single equilibrium that consists of applying to the same firm. Policies could cater to inducing forgetfulness or maintaining certain wage conditions that allow workers to learn to converge to the more efficient equilibria.

The second market structure consists of firms posting wages first, follow by workers making application choices. Workers adopt best response learning dynamics (Fudenberg and Levine (1998); Hopkins (1999)) by taking into account wages, beliefs about the other workers' actions and their bias towards firms based on their long-term experiences or perceived stereotypes. Given past realized actions, workers adjust their choice probabilities over the firms as their beliefs about each other evolve. Firms, tracking the changes in workers' beliefs, optimize wages in response to how they expect workers would behave in each period. In the long run, workers' and firms' behaviour converge to equilibria similar to the Quantal Response Equilibria (QRE). Multiple equilibria could exist when wage sensitivity and payoffs are high, and that bias from experiences are sufficiently weak. Workers are generally more likely to converge to asymmetric equilibria of applying with higher probability to different firms, but this can be affected by firms'

wage-setting behaviours and workers' sensitivity to expected payoffs. The asymmetric equilibria are also asymptotically stable if certain condition on wage dynamics is fulfilled, but strong bias from experiences may destabilize the system. In view of this, policies could tackle firms' wage-setting mechanism and workers' experience bias to induce more efficient and relatively stable outcome.

This paper shows that there could be labour market mismatch for substantial number of periods before workers converge to applying more efficiently to different firms, even when wages are dynamically changing. Workers' search strategies are affected by experiences, which can serve both as a catalyst and an inertia for reaching efficient outcome. The lack of switching in jobs observed in reality could happen when workers do not observe wages, and they naturally sort into different jobs based on their application experiences. On the other hand, when they observe wages, such lock-in effect may be less evident if they are sensitive to payoffs. However, when bias from long-term experiences are strong, workers' strategies could be resistant to changes even with wage information.

The rest of the paper is structured in the following manner: Section 2 introduces search with learning from feedback. Section 3 explores search with sequential learning. Section 4 conclude the paper by discussing the implications and extensions of the work.

Related literature. This work rides on top of extensive literature on experience-based learning to provide an evolutionary narrative for the labour market dynamics. There are many possible learning dynamics (see Fudenberg and Levine (1998)), but Erev and Roth (1998) shows that individuals display learning pattern close to reinforcement learning in experiments, thus this learning mechanism is adopted as the baseline in this paper to model job search behaviour. However, Hopkins (2002) highlights that when comparing between different learning models, this force of habit model is statistically insignificant in explaining Van Huyck et al. (1997)'s experimental results. Camerer and Hua Ho (1999) also postulates that an experience-weighted learning approach may fit individuals' learning pattern better. It would account for both actual and counterfactual payoffs if a worker select a different action. However, if workers cannot observe payoff from unchosen action, then such mechanism will not be feasible. Nonetheless, when workers can fully observe wages ex-ante to application decisions, there could be merit in accounting for opponent's empirical frequency of choices as one can infer the payoff of action not chosen. Therefore, I explore best response dynamics (Fudenberg and Levine (1998); Hopkins (1999)) that consider for expected payoffs given anticipated choice probabilities.

There are limited job search models that tracks transition path of individual firms' and workers' strategies based on application experiences. Studies have explored working experiences, which is not completely analogous to application experiences. For instance, Burdett and Mortensen (1998) proposes on-the-job experiences, which affect workers' preferences for firms, mainly due to human capital accumulation and expectation of higher wages as they climb the corporate ladder. Experiences in application stage, however, does not affect skills, yet it could also influence how one chooses firms, which constitute as possible means of directing search other than wages. Closer to application experiences is experiential job search highlighted by Kanfer and Bufton (2018), where workers adapt based on past involuntary job loss. In more direct relevance, Wanberg et al. (2020) examines past application experiences on one's adjustment of search behaviours. However, the work focuses more on empirical snapshots and descriptive analysis, an unified framework in modelling workers' adaptive learning process could be beneficial in formalizing labour market dynamics.

Furthermore, a core objective of integrating learning theory into job search is to offer some insights on equilibrium selection. Despite the presence of multiple equilibria, job search literature often focus on the symmetric equilibrium, where workers adopt the same application strategy, in both incomplete information (McCall (1970)) and complete information (Wright et al. (2021)), as

well as for some intermediary case of information availability (Wu (2020)). Although equilibria consist of workers applying to different firms are more efficient, even for partial information availability (Lu (2024)), these are often overlooked due to the perception that they are more difficult to coordinate on, which is a common approach in directed search literature (Galenianos and Kircher (2009)). However, experience-based learning suggests potential mechanisms for equilibrium selection, and provide some new insights on if workers could or could not converge to more efficient one-to-one matching outcome, and under what conditions.

This work also contributes to understanding some puzzles in real world observations. For instance, the phenomenon of workers' lack of job switch due to fundamental inertia in application decisions (Archer (2016)), as well as sorting behaviours displayed by different genders resulting from accumulation of experiences and beliefs built over the generations (Barbulescu and Bidwell (2013)), which differs from sorting due to skills (Eeckhout (2018)) or risk preferences (Fouarge et al. (2014)), all these may be supported by reinforcement learning behaviour, which exhibits familiarity-based learning pattern (Hopkins (2007)). Learning models could also substantiate the empirical evidence in Vafa et al. (2022), which shows that workers follow certain career trajectory and that past experiences are predictive of the jobs they end up in. While the predictive tool highlights observed job uptake pattern, it does not necessarily imply that workers do not attempt to apply to different jobs. Experience-based learning in the application stage formalize a possible channel behind such sorting behaviour. As a result, this could provide basis for policies to tackle potential mismatch due to sorting behaviour displayed in application choices.

Last but not least, job search models often concentrated on mapping from payoffs to choices, and the equilibrium choice probability distribution is often on the aggregate-level. Moen (1997) shows the probability distribution of workers choosing which sub-market to apply to. Workers have deterministic choices given payoffs, choices are probabilistic only at the population-level. While Galenianos and Kircher (2009) models workers' probabilistic choices, it places restrictions like symmetry in strategies to make analysis tractable, and is more concerned with mapping payoffs to aggregate choice distributions rather than to individual strategies. Therefore, this paper investigates specifically the individual-level choice probabilities rather than market-level aggregates, effectively showing impact of experiences on strategy formulations, which can be history and path-dependent. Moreover, investigating micro-level decision-making could help to inform policies influencing individual choices that can have macro implications.

2 Search with Learning from Feedback

In this section, I propose a market structure that conveys a traditional offline job search process, or search with online platforms but wage information are not explicitly revealed. In this setting, workers do not observe wages or strategies taken by other workers, who are simultaneously applying for jobs. They only learn the payoff from jobs they have applied to, and their search behaviours are solely affected by feedback received from previous periods. This learning process bear resemblance to Hopkins (2007), where consumers only receive payoff information of goods they have purchased. Using such learning mechanism, I investigate the impact of experiences on labour market dynamics.

2.1 Experience-driven Job Search with Fixed Wages

In a 2×2 set-up with 2 firms and 2 workers. Workers are indexed by $i = 1, 2$, and firms by $j = 1, 2$. Each firm offers one vacancy.

Timing and Set-up. Firms observe exogenous realization of productivity, denoted by $\mathbf{z} = \{z_1, z_2\}$, $\mathbf{z} \in Z$. They set wages, $\mathbf{w} = \{w_1, w_2\}$, which are fixed throughout all periods. Time is discrete, $t = \{0, 1, \dots, T\}$, representing generations of workers and firms.¹ In each period:

1. Workers do not observe the wages ex-ante to their selection between firm 1 and 2. But they can recall from past “self” or learn from the previous generation the application strategy, the realized choice and the payoffs. (For example, worker 1 of $t = 1$ learns from worker 1 of $t = 0$, same applies for worker 2.)
2. They devise an application strategy based on all the information.
3. Once workers made a choice, they observe a payoff, update their choice probabilities of each firm following Erev and Roth (1998)’s reinforcement learning algorithm. They then drop out of the market, never to return.

Firms’ side. In this fixed wage environment, wages are assumed to be rigid, such that firms set wages at the beginning of all application rounds and do not change them. Hereafter, suppose fixed wages satisfy $2w_1 > w_2 > \frac{w_1}{2}$, and that wages are bounded, $z_1 \geq w_1 \geq 0$, $z_2 \geq w_2 \geq 0$, such that wages are feasible and there can be multiple equilibria (as illustrated in the following workers’ side) for purpose of exploring equilibrium selection.

Workers’ side. Payoff structure faced by the workers is stationary, changes in workers’ choice probabilities over the firms are solely based on their reward observation from the previous period. (Nowé et al. (2012)) The search problem faced by the workers can be simply illustrated as a standard coordination game, where payoffs are fully revealed in the equilibrium. Assuming no obvious heterogeneity between workers, which could influence their probability of being hired, they would have equal probability of getting hired if applying to the same firm.

		Worker 2	
		F1	F2
Worker 1	F1	$\frac{w_1}{2}, \frac{w_1}{2}$	w_1, w_2
	F2	w_2, w_1	$\frac{w_2}{2}, \frac{w_2}{2}$

Table 1: Application Game Faced by Workers

¹For approximation of stochastic process in later parts, I consider $T \rightarrow \infty$, implying infinitely many generations.

The application game could consist of three Nash Equilibria (NEs): $(F1, F2)$ and $(F2, F1)$ if $2w_1 > w_2 > \frac{w_1}{2}$; and a mixed NE, $(F1, F2; \frac{2w_1-w_2}{w_1+w_2}, \frac{2w_2-w_1}{w_1+w_2})$ if $2w_1 \geq w_2 \geq \frac{w_1}{2}$.

Based on Duffy and Hopkins (2005)'s market entry game, as well as van Strien (2022) and Erev and Roth (1998), I define workers' actions, strategies, rewards, choice rule and update rule:

- **Actions.** For worker i , $A^i = \{F1, F2\}$, where $F1$ represents applying to firm 1 and $F2$ to firm 2. Same applies for worker $-i$, $A^{-i} = \{F1, F2\}$.
- **Strategies.** Worker i 's strategy is denoted as $\Delta_i = \{x_t = (x_{1t}, x_{2t})\}^T$, where $\sum_{j=1}^2 x_{jt} = 1$; and worker $-i$'s is $\Delta_{-i} = \{y_t = (y_{1t}, y_{2t})\}^T$, where $\sum_{j=1}^2 y_{jt} = 1$ (worker 2). x_{jt} and y_{jt} are the probabilities of firm j being chosen at time t . x_t is a pure strategy if $x_{jt} = 1$ for some j , similarly for y_t .
- **Rewards.** Wage matrix faced by worker i is $W = \begin{pmatrix} \frac{w_1}{2} & w_1 \\ w_2 & \frac{w_2}{2} \end{pmatrix}$. Payoff is denoted as $\pi_t^i = \pi_t^i(a_t^i, a_t^{-i})$, where $a_t^i \in A^i, a_t^{-i} \in A^{-i}$ are observed actions at the end of period t .

The reason behind this payoff structure is that I assume workers obtain positive reinforcement from both successful and unsuccessful application alike. This could arise from "good feelings" after being accepted for a job or just being interviewed, which is related to feeling validated and recognized professionally. (Briñol and Petty (2022)) When workers apply to the same firm, although one of them is not selected, they will still receive valuable information about their prospect of being hired. The positive signal of intent to hire is weighted by the probability of successful hire ($\frac{1}{2}$ for homogeneous workers), resulting in partial reinforcement of workers' choices if both applies to the same firm. Workers' choice in such case is updated based on potential payoff that encompasses the level of competition.

Workers' probability of selecting firm j at time t , x_{jt} and y_{jt} , depends on the propensities, denoted by q_{jt}^i . Each worker is endowed with an initial propensity for each action, $q_{j0}^i = \{q_{10}^i, q_{20}^i\}$. These initial propensities can be perceived as workers' respective innate preference before entering the labour market, and for subsequent periods, propensities can be referred to as accumulated payoffs obtained by selecting each firm (Beggs (2005)).

Herein, I adopt a linear **choice rule**:

$$\text{Worker } i: x_{jt} = \frac{q_{jt}^i}{\sum_{j=1}^J q_{jt}^i}, \text{ Worker } -i: y_{jt} = \frac{q_{jt}^{-i}}{\sum_{j=1}^J q_{jt}^{-i}} \quad (1)$$

Following Erev and Roth (1998) (ER) reinforcement learning mechanism, which specifies the **update rule** on how propensities are updated in each round:

$$\text{Worker } i: q_{j(t+1)}^i = q_{jt}^i + \pi_{jt}^i(a_t^i, a_t^{-i}), \text{ Worker } -i: q_{j(t+1)}^{-i} = q_{jt}^{-i} + \pi_{jt}^{-i}(a_t^i, a_t^{-i}) \quad (2)$$

When worker 1 select firm 1 in period t , if a positive feedback is received, then in the next period, the propensity of applying to firm 1 by worker 1 will increase by an increment equal to the realized payoff ($\pi_{1t}^i(F1, a_t^{-i})$) given observed actions ($a_t^i = F1, a_t^{-i}$).

For example, if worker 1 choose to apply to firm 1 in period 1 and worker 2 to firm 2, as payoffs (w_1, w_2) are revealed at the end of the period, action $F1$ and $F2$ are reinforced by the magnitude of w_1 and w_2 for worker 1 and 2, respectively. If both workers choose to apply to firm 1, action $F1$ will be reinforced by $\frac{w_1}{2}$ for both workers. Workers learnt the wage offered by firm 1 and how many candidates are competing for the firm. Even if they did not obtain the job, their propensity to firm 1 is positively affected because they gained some information, and if they

obtained the job, the fact that another candidate was also competing for the job implies higher possibility of not getting hired, so propensity to select the same firm is slightly discounted.

Workers only receive feedback from actions they actually take. The impact of the other worker's application strategy is implicit. One do not directly observe the choices taken by the other worker in the same period, which is realistic to the labour market. But worker $-i$'s choice in period t affects worker i in the same period via realized payoffs, thus influencing worker i 's choice propensities in period $t + 1$.

Given payoff matrix for worker 1 (W) and worker 2 (W^T):

$$W = \begin{pmatrix} \frac{w_1}{2} & w_1 \\ w_2 & \frac{w_2}{2} \end{pmatrix}, W^T = \begin{pmatrix} \frac{w_1}{2} & w_2 \\ w_1 & \frac{w_2}{2} \end{pmatrix} \quad (3)$$

I formulate the **expected change in application strategies** using Lemma 1 of Hopkins (2002) with the choice rule (1) and the update rule (2):

$$\text{Worker } i: E(x_{t+1}|q_t^i) - x_t = \frac{R(x_t)W y_t}{Q_t^i} + O\left(\frac{1}{(Q_t^i)^2}\right) \quad (4)$$

$$\text{Worker } -i: E(y_{t+1}|q_t^{-i}) - y_t = \frac{R(y_t)W^T x_t}{Q_t^{-i}} + O\left(\frac{1}{(Q_t^{-i})^2}\right) \quad (5)$$

where $q_t^i = \{q_{1t}^i, q_{2t}^i\}$, $q_t^{-i} = \{q_{1t}^{-i}, q_{2t}^{-i}\}$, each comprises of the propensities for the two actions at time t , corresponding to the two workers. $R(\cdot)$ is the replicator operator, reflecting how strategies evolve based on their relative payoffs:

$$R(x_t) = \begin{pmatrix} x_{1t}(1 - x_{1t}) & -x_{1t}x_{2t} \\ -x_{2t}x_{1t} & x_{2t}(1 - x_{2t}) \end{pmatrix}, R(y_t) = \begin{pmatrix} y_{1t}(1 - y_{1t}) & -y_{1t}y_{2t} \\ -y_{2t}y_{1t} & y_{2t}(1 - y_{2t}) \end{pmatrix}.$$

This implies, for instance, as worker 1's choice probability to firm 1 (x_{1t}) increases, its growth rate is dampened, and at the same time, the choice probability to firm 2 is also reduced. This determines how quickly the probability changes and which equilibrium the system moves towards. Lastly, $Q_t^i = \sum_{j=1}^J q_{jt}^i$ denotes the sum of propensities, it may be interpreted as a control over the magnitude of updates. As Q_t^i grows over time, the system would approximate continuous time dynamics in the limit. This relates to the step size, which describes the rate at which workers updates their strategies. In this model, workers would have different step size determined by their payoff experiences. But as t increases, with $Q_t^i \rightarrow \infty$ and $Q_t^{-i} \rightarrow \infty$, the effective step size is of order $\frac{1}{t}$, adjustments to strategies become less significant over time.

$$E(x_{t+1}|q_t^i) - x_t \approx \frac{R(x_t)W y_t}{Q_t^i}, E(x_{t+1}|q_t^i) - x_t \rightarrow 0 \quad (6)$$

$$E(y_{t+1}|q_t^{-i}) - y_t \approx \frac{R(y_t)W^T x_t}{Q_t^{-i}}, E(y_{t+1}|q_t^{-i}) - y_t \rightarrow 0 \quad (7)$$

Given the possibility of three equilibria in this setting, workers could either coordinate on applying to different firms or adopting a mixed strategy that suggest more randomized search behaviour. Based on Hopkins and Posch (2005), ER learning rule would not converge to a Nash Equilibrium (NE) linearly unstable under the replicator dynamics and it cannot converge to a rest point that is not a NE. The mixed strategy equilibrium is unstable under replicator dynamics, any perturbation would cause workers to drift towards one of the pure strategy equilibria. As a result, workers should eventually learn to coordinate on applying to different firms, and experiences would act as a natural selection mechanism for arriving at efficient outcome of one-to-one matching. However, the exact pure NE that is reached could be dependent on the initial conditions.

Proposition 1 (Stability of Asymmetric Equilibria). *Let (x_t, y_t) denote the stochastic process of workers' choice probabilities under reinforcement learning dynamics. The stochastic process can be approximated by the deterministic system in the limit of $t \rightarrow \infty$. Given fixed wages that satisfy $2w_1 > w_2 > \frac{w_1}{2}$, asymmetric equilibria are locally asymptotically stable.*

Proof. Based on equations (6) and (7), as $t \rightarrow \infty$, $Q_t^i \rightarrow \infty$ and $Q_t^{-i} \rightarrow \infty$. The expected change in strategies vanishes and the drift stabilizes, the stochastic process thus converges to the deterministic system:

$$\frac{dx_{1t}}{dt} = \lim_{Q_t^i \rightarrow \infty} Q_t^i (E(x_{t+1}|q_t^i) - x_t) = R(x_t)W y_t \quad (8)$$

$$\frac{dy_{1t}}{dt} = \lim_{Q_t^{-i} \rightarrow \infty} Q_t^{-i} (E(x_{t+1}|q_t^{-i}) - y_t) = R(y_t)W^T x_t \quad (9)$$

Given payoff matrices, W and W^T , from (3),

$$\frac{dx_{1t}}{dt} = f(x_{1t}, y_{1t}) = x_{1t}(1 - x_{1t})\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]y_{1t} + w_1 - \frac{w_2}{2} \quad (10)$$

$$\frac{dy_{1t}}{dt} = g(x_{1t}, y_{1t}) = y_{1t}(1 - y_{1t})\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]x_{1t} + w_1 - \frac{w_2}{2} \quad (11)$$

The Jacobian matrix:

$$J = \begin{bmatrix} \frac{\partial f}{\partial x_{1t}} & \frac{\partial f}{\partial y_{1t}} \\ \frac{\partial g}{\partial x_{1t}} & \frac{\partial g}{\partial y_{1t}} \end{bmatrix} = \begin{bmatrix} (1 - 2x_{1t})\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]y_{1t} + w_1 - \frac{w_2}{2} & x_{1t}(1 - x_{1t})\left(-\frac{w_1}{2} - \frac{w_2}{2}\right) \\ y_{1t}(1 - y_{1t})\left(-\frac{w_1}{2} - \frac{w_2}{2}\right) & (1 - 2y_{1t})\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]x_{1t} + w_1 - \frac{w_2}{2} \end{bmatrix} \quad (12)$$

For symmetric strategy, $x_{1t} = y_{1t}$, denote this as $\hat{x}y$, evaluating eigenvalues using $\det(J - \lambda I) = 0$:

$$\lambda = (1 - 2\hat{x}y)\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]\hat{x}y + w_1 - \frac{w_2}{2} \pm \sqrt{(\hat{x}y(1 - \hat{x}y)\left(-\frac{w_1}{2} - \frac{w_2}{2}\right))} \quad (13)$$

Given $2w_1 > w_2 > \frac{w_1}{2}$ is satisfied, there will be a positive and a negative eigenvalue, making this a saddle point. In the special case of homogeneous firms, $z_1 = z_2 = z$, $w_1 = w_2 = w^*$, at $x_{1t} = y_{1t} = \frac{1}{2}$:

$$J\left(\frac{1}{2}, \frac{1}{2}\right) = \begin{bmatrix} 0 & -0.25w^* \\ -0.25w^* & 0 \end{bmatrix} \quad (14)$$

$\lambda = \pm 0.25w^*$. This is a saddle point and is unstable for any positive wages.

For asymmetric strategies, $x_{1t} = 1, y_{1t} = 0$ or $x_{1t} = 0, y_{1t} = 1$, the eigenvalues are:

$$\lambda_1 = (1 - 2x_{1t})\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]y_{1t} + w_1 - \frac{w_2}{2}, \lambda_2 = (1 - 2y_{1t})\left[-\frac{w_1}{2} - \frac{w_2}{2}\right]x_{1t} + w_1 - \frac{w_2}{2} \quad (15)$$

for $2w_1 > w_2 > \frac{w_1}{2}$, at $x_{1t} = 1, y_{1t} = 0$, $\lambda_1 = -w_1 + \frac{w_2}{2} < 0$, $\lambda_2 = \frac{w_1}{2} - w_2 < 0$; at $x_{1t} = 0, y_{1t} = 1$, $\lambda_1 = \frac{w_1}{2} - w_2 < 0$, $\lambda_2 = -w_1 + \frac{w_2}{2} < 0$. These points are locally asymptotically stable. \square

Example 1. (Unstable Mixed Strategy Equilibrium) Suppose $w_1 = 3$, $w_2 = 3$, there could exist 3 possible equilibria: $(F1, F2)$, $(F2, F1)$, and $(F1, F2, \frac{4}{5}, \frac{1}{5})$. Assuming $Q_t^i = Q_t^{-i} = 1$,

$$\frac{dx_{1t}}{dt} = f(x_{1t}, y_{1t}) = x_{1t}(1 - x_{1t})(2 - 2.5y_{1t}) \quad (16)$$

$$\frac{dy_{1t}}{dt} = g(x_{1t}, y_{1t}) = y_{1t}(1 - y_{1t})(2 - 2.5x_{1t}) \quad (17)$$

$$J(x_{1t}, y_{1t}) = \begin{bmatrix} 0 & -0.4 \\ -0.4 & 0 \end{bmatrix} \bigg|_{\frac{4}{5}, \frac{4}{5}} \quad (18)$$

Evaluating at the point of $(x_{1t}, y_{1t}) = (\frac{4}{5}, \frac{4}{5})$, $\lambda = \pm 0.4$.

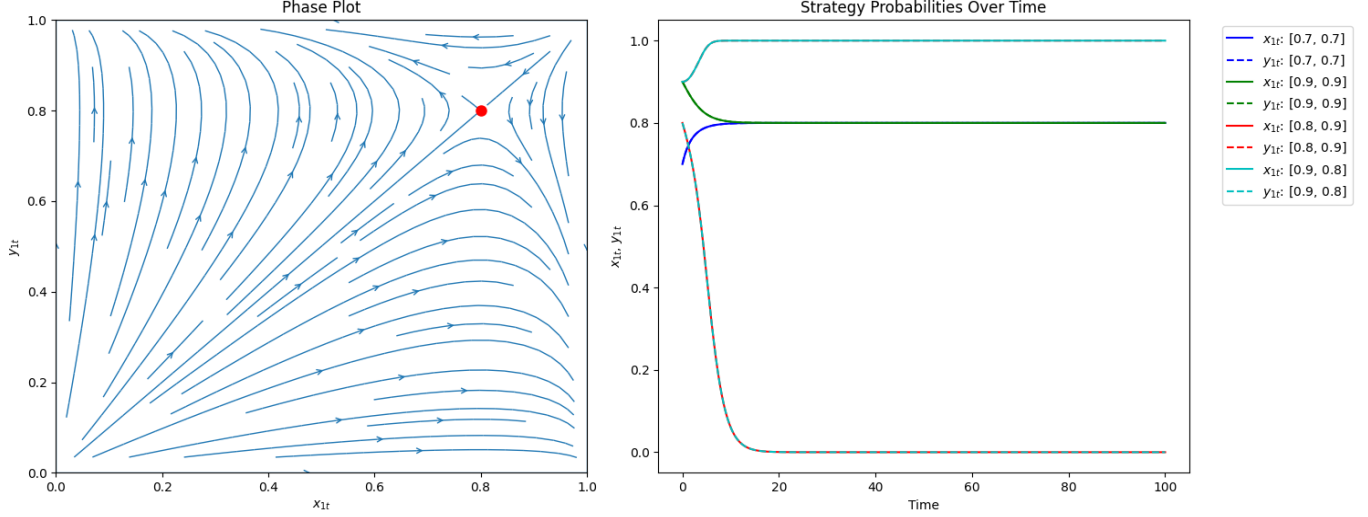


Figure 1: Illustration of Example 1

Figure 1 shows the phase plot and strategy change over time. Most learning trajectories demonstrate workers should end up in one of the pure strategy equilibria.

Erev and Roth (1998) highlighted one special feature of this learning mechanism – its heavy reliance on initial propensities and decreasing impact from recent experiences as accumulated payoffs become higher.

Observation 1 (Initial Bias and Experiences on Equilibrium selection). *Equilibrium selection is influenced by initial bias, characterized by initial propensities (q_{j0}^i, q_{j0}^{-i}) ; as well as initial experiences, characterized by feedback received in the first few application rounds:*

- *Strong prior preference or goodwill towards specific firm could create fundamental inertia for workers to adjust their application strategies, contributing to immediate sorting into different firms or persistent overcrowding at a single firm.*
- *Workers could be locked-in to choices they made initially, leading to experience-based sorting.*

Given wage condition holds for multiple equilibria, $2w_1 > w_2 > \frac{w_1}{2}$, as $t \rightarrow \infty$, workers would eventually coordinate on applying to different firms.

$$\lim_{t \rightarrow \infty} x_{1t} = \bar{x} = 1 \text{ if } \lim_{q_1 \rightarrow \infty} \frac{q_1}{q_1 + q_2} = 1 \quad (19)$$

$$\lim_{t \rightarrow \infty} x_{1t} = \bar{x} = 0 \text{ if } \lim_{q_2 \rightarrow \infty} \frac{q_1}{q_1 + q_2} = 0 \quad (20)$$

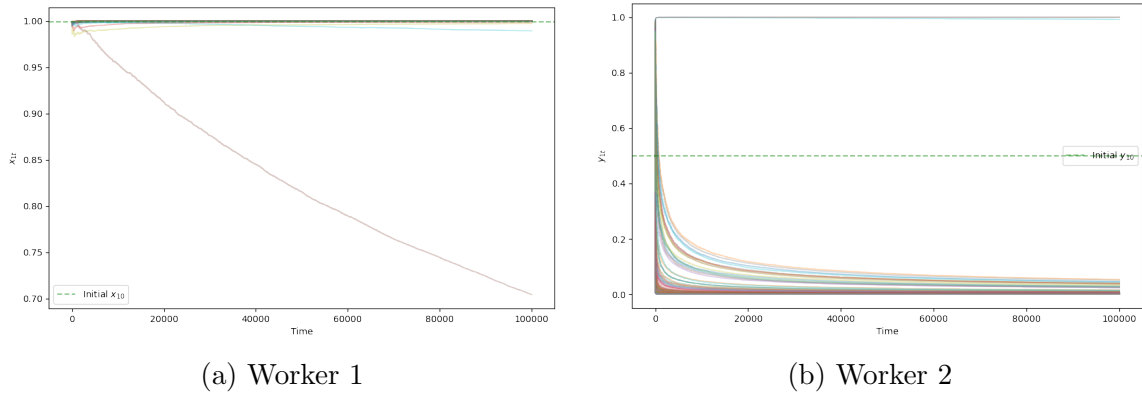
Despite this, there could be multiple periods where workers oscillate between choosing firm 1 and 2, leading to possible overcrowding for many periods over the trajectory of change in choice probabilities. Strong initial bias by both workers (i.e. high initial propensities towards the same firm) could also lead to biased search for substantial number of periods.

Hopkins (2007) also demonstrates that consumers could be locked into choosing certain good that they initially prefer and undergo familiarity-based learning. In the job search context, workers choosing to apply to firm 1 and 2 respectively in the initial periods would make $(F1, F2)$ more likely to be selected than other equilibria in the long run. This suggests that for homogeneous workers, who could start with random search, their diverse application strategies may be a result of randomness and luck, driven by positive reinforcement in initial periods of their job search. This also demonstrates that sorting behaviour can be experience-driven rather than based solely on skills.

Example 2. (*Heterogeneous Initial Bias*) Suppose initialization propensities are $q_{10}^i = 1000, q_{10}^{-i} = 1, q_{20}^i = q_{20}^{-i} = 1$, worker 1 denotes higher propensity to firm 1 than firm 2. Workers' choice probabilities to firm 1:

$$x_{10} = \frac{1000}{1001} \approx 0.999, y_{10} = \frac{1}{2} \quad (21)$$

Worker 1 applies with higher probability to firm 1 by default as compared to worker 2. As t increases, $(F1, F2)$ is more likely to be reached.



Figures show workers' application probability to firm 1 (y-axis) against number of periods (x-axis) for 10 simulation sessions ($w_1 = w_2 = 5, t = 100000, q_{10}^i = 1000, q_{20}^i = q_{10}^{-i} = q_{20}^{-i} = 1$).

Figure 2: Learning Path of Worker 1 and 2.

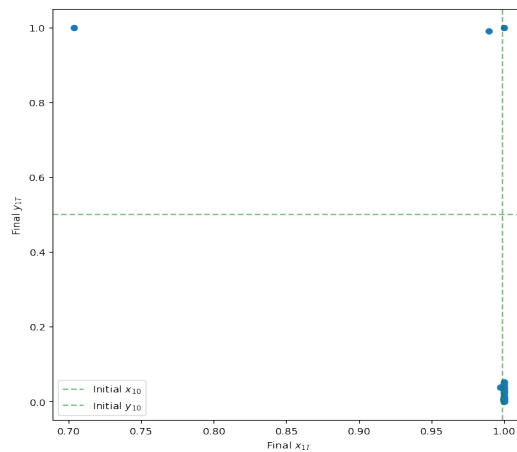


Figure 3: Workers' Final Choice Probability in Period T

Figure 2a & 2b demonstrate simulations of how workers' choice probabilities of firms evolve over time given past experiences. Worker heterogeneity, on the basis of differences in initial

propensities, could lead to sorting behaviour. In this example, it is more likely for worker 1 to choose firm 1 during initial rounds, thus more likely to be lock-in to firm 1. Since oscillations are fundamental characteristic of this learning mechanism, there is always a positive chance of selecting an alternative action, thus it remain probable for $(F2, F1)$ to be selected in the long run, but the likelihood of converging to $(F1, F2)$ is higher. Figure 3 shows the final choice probability in period T , most of the sessions stop at close to $(F1, F2)$.

In Barbulescu and Bidwell (2013), students with same education backgrounds could display segregation in job applications due to gender stereotypes associated with the jobs; and Terjesen et al. (2007) also found that females, in comparison to male counterparts in universities, put greater weights on “using your degree skills”, thus are more likely to go for jobs related to their degree. In the learning model, these can be interpreted as worker heterogeneity in initial propensities, which translate into dispersion in choice probabilities and affect equilibrium selection.

2.1.1 Partial Recall of Experiences

Perfect recall of experiences may be a stringent assumption, individuals are often subjected to limited memory. Cognitive load theory suggests individuals’ working memory is constrained to a capacity of approximately 4 elements of information (Paas and Ayres (2014)). To implement the impact of partial recall on workers, I include a forgetting parameter, η , $\eta \in (0, 1)$, to account for recency effect. As a result, past experiences or knowledge could have a diminishing effect on current application decisions (Erev and Roth (1998)).

$$q_{j(t+1)}^i = (1 - \eta)q_{jt}^i + \pi_{jt}^i(a_t^i, a_t^{-i}) \quad (22)$$

Proposition 2 (Partial Recall). *Let $\eta \in (0, 1)$ be the experience decay parameter in the propensity updating process, for equation (22):*

1. *As $\eta \rightarrow 0$, the propensity updating process converges to perfect recall.*
2. *As $\eta \rightarrow 1$, next period propensity $q_{j(t+1)}^i$ converges to being determined solely by realized payoffs in t , $\pi_{jt}^i(a_t^i, a_t^{-i})$.*
3. *For $0 < \eta < 1$, weight placed on previous period propensities is lower, the influence of initial propensity q_{j0}^i diminishes as $T \rightarrow \infty$.*

Proof. For T periods, as $\eta \rightarrow 0$,

$$\lim_{\eta \rightarrow 0} q_{jT}^i = q_{j0} + \pi_{j0} + \pi_{j1} + \dots + \pi_{jT-2} + \pi_{j(T-1)} \quad (23)$$

As $\eta \rightarrow 1$,

$$\lim_{\eta \rightarrow 1} q_{jT}^i = \pi_{j(T-1)} \quad (24)$$

For $0 < \eta < 1$,

$$q_{jT}^i = (1 - \eta)^T q_{j0} + (1 - \eta)^{T-1} \pi_{j0} + (1 - \eta)^{T-2} \pi_{j1} + \dots + (1 - \eta) \pi_{jT-2} + \pi_{j(T-1)} \quad (25)$$

Coefficient on q_{j0}^i is $(1 - \eta)^T$, $\lim_{T \rightarrow \infty} (1 - \eta)^T = 0$. \square

In the extreme case of complete forgetting ($\eta = 1$), the probability of worker 1 selecting firm 1 depends solely on the last period feedback:

$$\lim_{\eta \rightarrow 1} x_{1T} = \frac{\pi_{1(T-1)}^i}{\pi_{1(T-1)}^i + \pi_{2(T-1)}^i}, \lim_{\eta \rightarrow 1} y_{1T} = \frac{\pi_{1(T-1)}^{-i}}{\pi_{1(T-1)}^{-i} + \pi_{2(T-1)}^{-i}} \quad (26)$$

Since workers can only select one firm in specific period, this implies for worker 1, either $\pi_{1(T-1)}$ or $\pi_{2(T-1)}$ will be positive, therefore, either $x_{1T} \rightarrow 1$ or $x_{1T} \rightarrow 0$. There is convergence to applying with certainty to firm 1 or 2. Same applies for worker 2.

For intermediary forgetfulness ($0 < \eta < 1$), past events will have some impact on application choices. But as compared to perfect recall, past events have diminishing impact on propensities, which could lead to slower augmentation of Q_t^i and Q_t^{-i} . Therefore, it may take longer for the system to settle into equilibrium. Whilst in the fixed wage environment, workers are expected to converge to pure NEs, limited memory could contribute to longer process of reaching efficient outcomes and there may be more instances of mismatch.

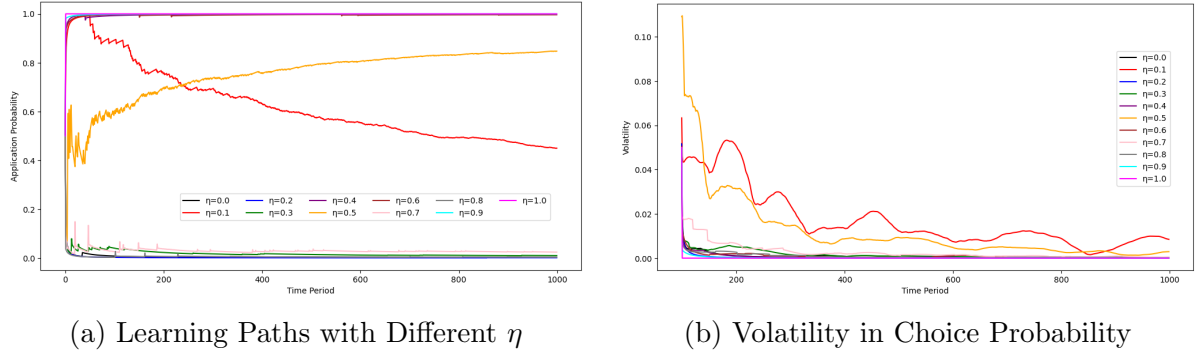


Figure shows worker 1's application probability and volatility in choice probability for different η , $\eta \in [0, 1]$ across first 1000 periods. ($w_1 = w_2 = 5$, $t = 100000$, $q_{10}^i = q_{20}^i = q_{10}^{-i} = q_{20}^{-i} = 1$).

Figure 4: Worker 1's Choice Probability of Firm 1 with Forgetting

In Figure 4a, I show a simulation of worker 1's choice probability of firm 1 for different forgetting parameters. When workers do not remember any past events ($\eta = 1$), choice probability go straight to 1 or 0. It is possible that if both workers apply to firm 1 and receive positive feedback of $\frac{w_1}{2}$ end up overcrowding at firm 1. When workers retain some memories ($\eta < 1$), there is convergence to pure NEs, but choice probabilities are noisy due to positive weight placed on past propensities. For instance, Figure 4b shows high volatility² in choice probability when $\eta \neq 1$. While intuitively, larger forgetting parameter should imply more volatile choice probabilities, but in this setting, since propensities depend on realized payoffs, which are inherently stochastic when choice probabilities are not deterministic, so larger forgetting parameter may not imply more volatile outcomes.

2.1.2 Policy Implications

In a static wage environment, where wages satisfy the condition for multiple equilibria to exist, workers would eventually coordinate on applying to different firms in the long run. Experiences could serve as a natural mechanism paving towards an efficient market outcome. However, there could be many periods of mismatch before pure NEs are reached, thus suggesting policies could tackle the learning trajectory.

The convergence can be slower when workers have strong initial bias to the same firm, indicating policies could cater to reducing initial bias to facilitate faster coordination. Eliciting different prior preferences also have some merits as having heterogeneous initial bias could be beneficial to facilitate faster coordination on applying to different firms.

²Higher volatility implies how much application probability to firm 1 in t is higher than average choice probability of the past 100 periods (between t and $t - 99$).

Furthermore, in this learning model, initial experiences are important for later choices and could influence the equilibrium selected. Once workers start with a job, then they are likely to settle into similar jobs in the future. If both workers apply for the same job in the initial periods and obtain similar experiences, it will take substantial number of periods for them to adjust, policies could speed up the process by lowering the memory impact. This may be done by influencing the degree of forgetting through platform recommendation system that control how much exposure one has to past information when searching for jobs. However, one need to determine whether to induce more or less forgetfulness. For instance, when workers have sufficiently diverse experiences, perfect recall could encourage faster convergence towards efficient outcome; but when workers' initial propensity are high and biased towards the same firm, or experiences from first few applications reinforce potential overcrowding, then there is benefits to "unlearn".

2.2 Two-sided Simultaneous Learning with Dynamic Wages

In this section, I relax the assumption on wage rigidity, such that firms can also be adaptive learners and wages are affected by workers' search behaviours. I explore how workers react to a dynamic wage environment, and if they would learn to coordinate on applying to different firms. I will also investigate how wages evolve.

In this set-up, both firms and workers are learning simultaneously and are assumed to adopt Erev and Roth (1998) reinforcement learning algorithm.

Workers' side. Workers follow the same learning pattern as Section 2.1.

Firms' side. Given exogenous realization of productivities, $\mathbf{z} = \{z_1, z_2\}$, firms select wages, $\mathbf{w} = \{w_1, w_2\}$, where $0 \leq \mathbf{w} \leq \mathbf{z}$. It is assumed that as a new firm in the labour market, it is unlikely for it to know which wage to set to attract workers at the beginning of application rounds, but over time, it would learn to choose wages based on the responses it obtained. For example, in time t , if firm 1 choose w_1 and receives both workers' application and is able to produce this period, then it is likely to select the same wage again in period $t + 1$ given an update rule. However, if it receives no worker in period t , then the choice of wage value w_1 is not reinforced. I define firms' actions, strategies, rewards, choice rule and update rule below:

- **Actions.** Firm j has finite and discrete number of actions, $A^j = (0, 1, 2, \dots, z_j)$. Assuming firm homogeneity, $A^j = A^{-j} = (0, 1, 2, \dots, z)$. Each action effectively corresponds to the wage offered, $a^j = w_j$.
- **Strategies.** Firm j 's strategy is denoted as $\Delta_j = \omega_t^j = (\omega_{0t}^j, \omega_{1t}^j, \dots, \omega_{zt}^j)^T$, where $\sum_{a^j=0}^z \omega_{a^j t}^j = 1$; and firm $-j$'s is $\Delta_{-j} = \omega_t^{-j} = (\omega_{0t}^{-j}, \omega_{1t}^{-j}, \dots, \omega_{zt}^{-j})^T$, where $\sum_{a^j=0}^z \omega_{a^j t}^{-j} = 1$. $\omega_{a^j t}^j$ and $\omega_{a^j t}^{-j}$ are the probabilities of each action a^j being chosen by firm 1 and 2 at time t .
- **Rewards.** If at least one worker applies to firm j , firm j receives a payoff of $z_j - w_j$, otherwise 0. Two possible outcomes:

$$I_i^j = \begin{cases} 1 & \text{if worker } i \text{ chooses firm } j \\ 0 & \text{otherwise} \end{cases}, \quad I_{-i}^j = \begin{cases} 1 & \text{if worker } -i \text{ chooses firm } j \\ 0 & \text{otherwise} \end{cases}$$

- **Choice Rule.** Firms' choice probabilities of selecting each action at time t , $\omega_{a_t}^j$ and $\omega_{a_t}^{-j}$, depend on propensities of each action being chosen, denoted as $\theta_{a_t}^j$ and $\theta_{a_t}^{-j}$ respectively.

Adopting a linear choice rule:

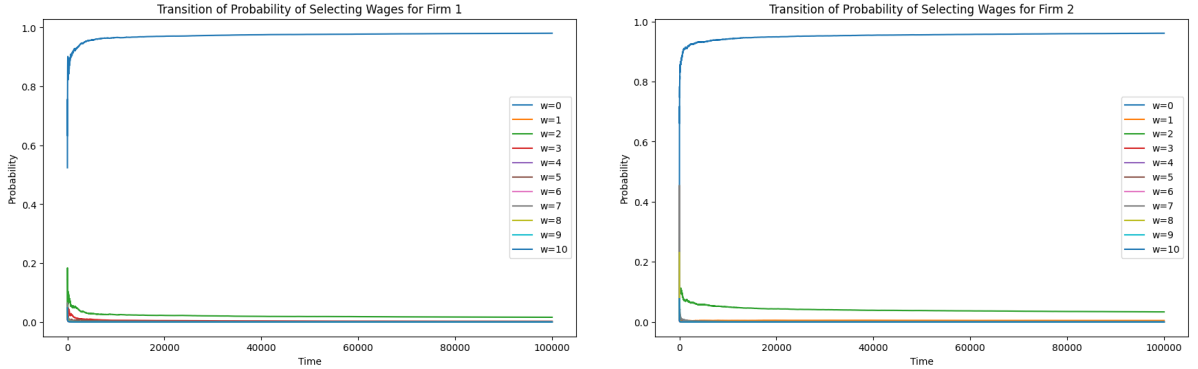
$$\text{Firm } j: \omega_{a_t^j}^j = \frac{\theta_{a_t^j}^j}{\sum_{a^j=0}^{z_j} \theta_{a_t^j}^j}, \text{Firm } -j: \omega_{a_t^{-j}}^{-j} = \frac{\theta_{a_t^{-j}}^{-j}}{\sum_{a^{-j}=0}^{z_{-j}} \theta_{a_t^{-j}}^{-j}} \quad (27)$$

- **Update Rule.** Propensities are updated based on realized payoffs:

$$\text{Firm } j: \theta_{a_{(t+1)}^j}^j = \theta_{a_t^j}^j + \pi_t^j(a_t^i, a_t^{-i}, a_t^j, a_t^{-j}), \text{Firm } -j: \theta_{a_{(t+1)}^{-j}}^{-j} = \theta_{a_t^{-j}}^{-j} + \pi_t^{-j}(a_t^i, a_t^{-i}, a_t^j, a_t^{-j}) \quad (28)$$

While workers gain positive reinforcement from both successful and unsuccessful application alike. For the firms, they only receive reinforcement for the wage they set when they successfully hire a worker. This is because they do not gain any additional information about how close a worker was to choosing them if they receive no application. As a result, workers and firms have slightly different learning dynamics due to differences in their access to information.

Suppose both firms and workers are homogeneous, they start with uniform probability over their action space. Figure 5 demonstrates growing probability of wage 0 being chosen by both firms. Figure 6 shows that in the final period of 10 simulated sessions (runs), the probabilities of setting each wage are higher for low wage values. Both figures suggest that wages will eventually be pushed down towards 0.



(a) Firm 1's Choice Probabilities of Wages

(b) Firm 2's Choice Probabilities of Wages

Figure shows firms' probability of choosing each discrete wage value, $\mathbf{w} \in [0, 10]$ for $t = 100000$, $q_{j0}^i = q_{j0}^{-i} = 1$,

$$\theta_{a_0^j}^j = \theta_{a_0^{-j}}^{-j} = 1, z_1 = z_2 = z = 10.$$

Figure 5: Learning Path of Firm 1 and 2 in 2-sided RL

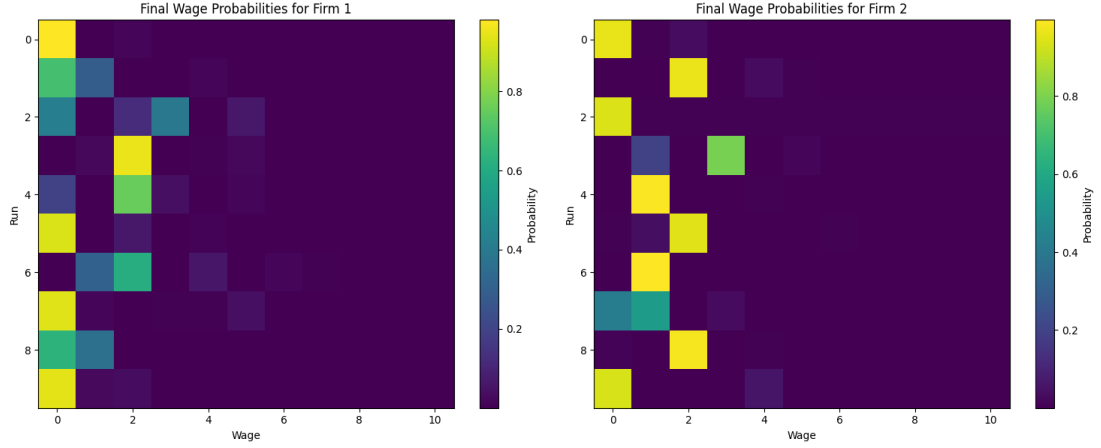


Figure 6: Firms' Probability of Selecting Each Wage in Period T Across 10 Sessions

For workers' side, Figure 7 shows the choice probabilities in the final period of 10 simulated sessions. There is higher likelihood of $(F1, F2)$ being selected. However, when compare to the static wage environment in Figure 3, the results are less saturated around pure NEs. There are also higher chances of workers applying to the same firm after substantial number of periods at end of period T , which are inefficient outcomes.

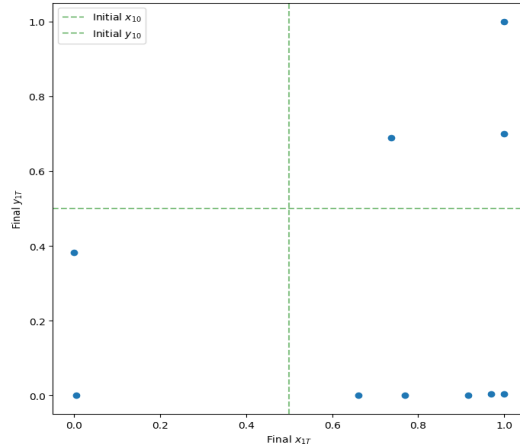


Figure 7: Workers' Final Choice Probability in Period T

Proposition 3 (Convergence in Dynamic Wage Environment). *Given workers' strategies, x_t and y_t , and firms' strategies, $\omega_{a_t^j}^j$ and $\omega_{a_t^{-j}}^{-j}$, at time t . The system converges to a fixed point, $(x^*, y^*, \omega_{a_t^j}^{j*}, \omega_{a_t^{-j}}^{-j*})$, as $t \rightarrow \infty$.*

Proof. For workers, I show in Section 2.1 that as $t \rightarrow \infty$, by equations 6 and 7, $E(x_{t+1}|q_t^i) - x_t \rightarrow 0$ and $E(y_{t+1}|q_t^{-i}) - y_t \rightarrow 0$.

For firms, their expected change in strategies are:

$$\text{Firm } j: E[\omega_{a_{t+1}^j}^j | \theta_{a_t^j}^j] - \omega_{a_t^j}^j = \frac{[1 - (1 - x_{jt})(1 - y_{jt})](z_j - a_t^j)}{S_t^j} + O\left(\frac{1}{S_t^{j2}}\right) \quad (29)$$

$$\text{Firm } -j: E[\omega_{a_{t+1}^{-j}}^{-j} | \theta_{a_t^{-j}}^{-j}] - \omega_{a_t^{-j}}^{-j} = \frac{[1 - (1 - x_{(-j)t})(1 - y_{(-j)t})](z_{-j} - a_t^{-j})}{S_t^{-j}} + O\left(\frac{1}{S_t^{-j2}}\right) \quad (30)$$

where $[1 - (1 - x_{jt})(1 - y_{jt})]$ reflects the probability that at least one worker applies to the firm; $S_t^j = \sum_{a^j=0}^{z_j} \theta_{a_t^j}^j$, $S_t^{-j} = \sum_{a^{-j}=0}^{z_{-j}} \theta_{a_t^{-j}}^{-j}$ are the sum of propensities over different actions in period t for each firm.

As $t \rightarrow \infty$, S_t^j and S_t^{-j} grow larger. From equations 29 and 30,

$$E[\omega_{a_{t+1}^j}^j | \theta_{a_t^j}^j] - \omega_{a_t^j}^j \approx \frac{[1 - (1 - x_{jt})(1 - y_{jt})](z_j - a_t^j)}{S_t^j}, E[\omega_{a_{t+1}^j}^j | \theta_{a_t^j}^j] - \omega_{a_t^j}^j \rightarrow 0 \quad (31)$$

$$E[\omega_{a_{t+1}^{-j}}^{-j} | \theta_{a_t^{-j}}^{-j}] - \omega_{a_t^{-j}}^{-j} \approx \frac{[1 - (1 - x_{(-j)t})(1 - y_{(-j)t})](z_{-j} - a_t^{-j})}{S_t^{-j}}, E[\omega_{a_{t+1}^{-j}}^{-j} | \theta_{a_t^{-j}}^{-j}] - \omega_{a_t^{-j}}^{-j} \rightarrow 0 \quad (32)$$

Both workers' and firms' expected change in strategies converge to 0. \square

As workers' strategies stabilize, wages will be driven down to 0 as positive wages reduce profits. At the point where wages are exactly 0, there will be a continuum of equilibria. Both Jacobian and Hessian matrix evaluated at the equilibrium points become 0, which means all equilibrium points are neutrally stable and are weak NEs. Based on this logic, the equilibrium selected depends on the direction of strategy adjustment at the time when wages are still positive. Once wages hit 0, strategies are "locked-in". Such outcome is seemingly "uninteresting" as both sides simply stop learning because there are no rewards and no adaptive pressure, this does not imply strategies are optimal. Despite the outcome, it may be interesting to explore the learning trajectory since equilibrium selection is path-dependent, efficient outcomes of one-to-one matching may be reached depending on the learning dynamics.

Along the learning path, since wages are dynamically changing, workers could be facing different "games" depending on the wage conditions. There could be equilibrium switching as one moves from one "game" to another.

Definition 2.1. (*Equilibrium Switching*) Consider the game faced by workers at time t to be G_t , the set of Nash Equilibria (NEs) associated with the specific game is defined to be $NE(G_t)$. Equilibrium switching is defined to occur if:

1. There exist two distinct equilibrium sets NE_1 and NE_2 , such that $NE(G_t) = NE_1$ for $t \leq \hat{t}$; $NE(G_t) = NE_2$ for $\hat{t} < t \leq T$, \hat{t} is the critical time of game change, T being total number of periods.
2. There can be multiple equilibrium switching over workers' learning trajectory.

Workers could effectively be facing three possible games for some given wage conditions (see Figure 8, 9, 10). They encompass different sets of NE(s). Equilibrium switching (Definition 2.1) could occur as payoffs evolve, and workers are learning different sets of NE(s) along the learning trajectory.

		Worker 2	
		F1	F2
Worker 1	F1	$\frac{w_1}{2}, \frac{w_1}{2}$	w_1, w_2
	F2	w_2, w_1	$\frac{w_2}{2}, \frac{w_2}{2}$

Figure 8: G1: $w_{1t} > 2w_{2t}$

		Worker 2	
		F1	F2
Worker 1	F1	$\frac{w_1}{2}, \frac{w_1}{2}$	w_1, w_2
	F2	w_2, w_1	$\frac{w_2}{2}, \frac{w_2}{2}$

Figure 9: G2: $2w_{1t} > w_{2t} > \frac{w_{1t}}{2}$

		Worker 2	
		F1	F2
Worker 1	F1	$\frac{w_1}{2}, \frac{w_1}{2}$	w_1, w_2
	F2	w_2, w_1	$\frac{w_2}{2}, \frac{w_2}{2}$

Figure 10: G3: $w_{2t} > 2w_{1t}$

Based on the analysis for static wage environment in Section 2.1, it is expected that if wage condition for $G2$ can be sustained for a long period of time, workers would be able to learn to converge to the pure NEs, leading to efficient outcomes. However, since one may have to overcome a large inertia due to accumulated propensities from previous game play, which could involve learning a different set of NE(s), therefore, even if equilibrium switching happens (e.g. from Figure 8 to 9), workers' actual behaviour from converging to one set of NE(s) to another happens with a time lag.

Proposition 4. (*Delayed Adaptation in Equilibrium Switching*) *Following equilibrium switching at time \tilde{t} due to evolving wages, workers will not immediately transition from learning NE_1 to NE_2 . There exists a time lag in adapting their strategies:*

For $t \leq \tilde{t}$, workers converge to NE_1 :

$$\lim_{t \rightarrow \infty, t \leq \tilde{t}} x_t, y_t \in \text{Support of } NE_1 \quad (33)$$

At $t = \tilde{t}$, game changes with NE_2 becomes the set of new equilibria.

For $\tilde{t} < t < \bar{t}$, workers can be described to be in the transitional state that are still influenced by NE_1 , but beginning to adapt to NE_2 :

$$x_t, y_t \in \text{Support of } NE_1 \cup NE_2 \quad (34)$$

For $t \geq \bar{t}$, workers fully converge to NE_2 :

$$\lim_{\bar{t} \rightarrow \infty, t \geq \bar{t}} x_t, y_t \in \text{Support of } NE_2 \quad (35)$$

where x_t, y_t are set of choice probability over firms in time t for worker 1 and 2, respectively.

Proof. Shown by Example 3. □

Given perfect memory, time lag is infinite, the system never fully adapt to new conditions due to accumulated past payoffs. In Figure 11, I show for a single simulation session with 100 periods, there is switching between different wage regimes, corresponding to the different games ($G1$, $G2$, $G3$).

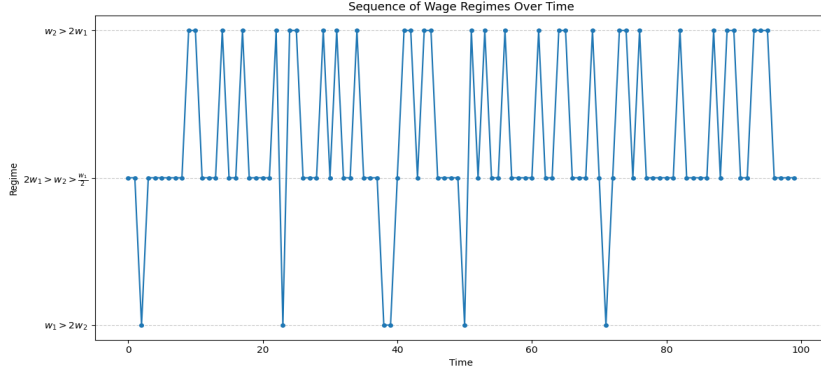


Figure shows wage transition across different regimes, characterized by wage conditions for ($z_1 = z_2 = 10$, $q_{10}^i = q_{20}^i = q_{10}^{-i} = q_{20}^{-i} = 1$), $t = 100$.

Figure 11: Switching Between Wage Regimes

In Figure 12, I illustrate the evolution of wages and choice probabilities. Suppose workers and firms are in regime 1 ($G1$) (red region), workers are converging towards $(F1, F1)$, firm 1's wage setting behaviour will be reinforced, and lower w_1 value will receive stronger reinforcement as profit ($z_1 - w_1$) received is higher. Decreasing w_1 could lead to regime switch. If a switch to regime 3 ($G3$) (green region) ensues, as workers learn to play $(F2, F2)$, lower w_2 will be reinforced more strongly, prompting another possible regime switch to $G2$. While workers' choice probabilities stabilizing in $G2$ (blue region) potentially prompt less incentive for further regime switch, stochasticity in workers' realized actions and equal attractiveness of $(F1, F2)$ and $(F2, F1)$ may lead to same action, such as both choosing firm 1, being realized, which could affect wages and incentivize regime switch. As a result, there may be constant regime change until workers' behaviour stabilizes and wages are driven down to 0.

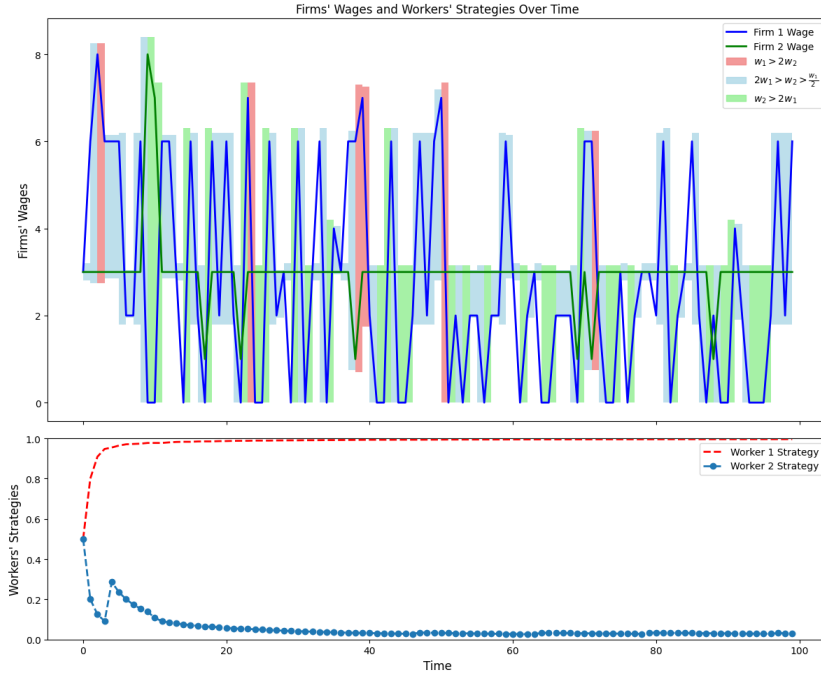


Figure 12: Changes in Wages and Workers' Choice Probabilities Over Time

Suppose I run the session for 10000 periods, and compute the conditional probability of switching

from one regime to another based on the empirical frequency:

$$P(G_{t+1} = G_i | G_t = G_j) = \frac{n_{ij}}{n_i} \quad (36)$$

where G_t refers to the game played in period t , reflective of the regime workers are in; n_{ij} is the counts of regime change from i to j ($i, j \in \{1, 2, 3\}$), and n_i is the counts of being in regime i . Figure 13 shows there is higher probability of switching from any regime to the G_2 regime for this simulation. However, given the stochasticity of action realizations by the workers and the firms due to their probabilistic behaviours, the transition patterns could vary across different simulation sessions.

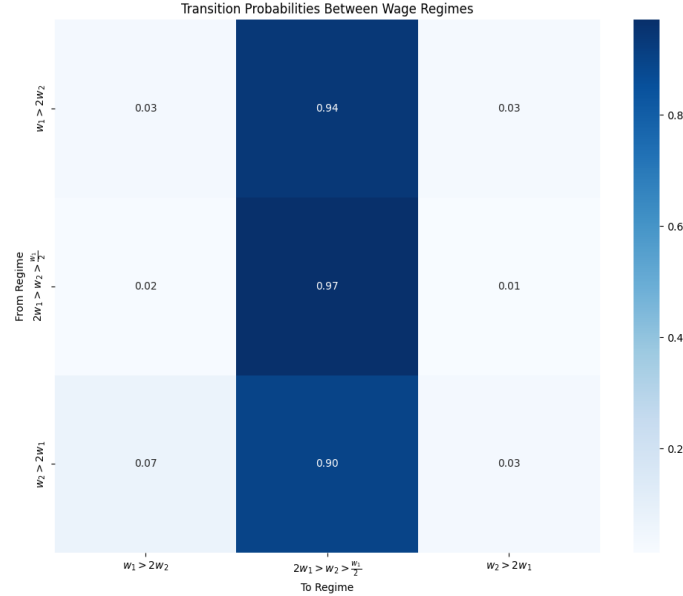


Figure 13: Transition Probability from Regime to Regime

Similar to static wage environment, dynamic wage environment could lead to prolonged periods of mismatch as workers need time to learn the NE(s). However, given wages vary over time, workers could be learning different sets of NE(s) as payoff structures shift over time. This could further contribute to mismatch as workers may need to overcome inertia from previous accumulated experiences as they adapt their search strategies amidst the new wage condition.

Another important question is equilibrium selection in the dynamic wage environment. Given that wages are pushed down to 0 in the long run, workers will be indifferent between selecting firm 1 or 2, and there will be a continuum of equilibria. Equilibrium selection is essentially path-dependent, and relies on what game is being played before wages hit 0. For one-to-one matching (i.e. $(F1, F2)$ or $(F2, F1)$) to be more likely, wage condition for $G2$ (Figure 9) needs to be maintained for substantial periods of time, such that pure NEs are possible outcomes where workers can learn to coordinate on.

2.2.1 Partial Recall of Experiences

Workers may not have perfect recall of past experiences, but firms are likely able to keep track of all past feedback by storing hiring information in a database that can be easily maintained and retrieved, and human resources tend to keep a record of wages offered and accepted. Therefore, I impose a memory decay factor on workers' side akin to Section 2.1.1 (Equation 22).

Since workers do not have perfect recall of past experiences, it is expected there will be slower convergence towards an equilibrium point. While longer learning trajectory could generate more instances of mismatch, the ability to forget “misaligned” market states as payoff structure shifts may also be an asset. As regime changes, the time lag for workers to start adapting to learn new set of NE(s) would be shorter under partial recall than perfect memory. However, discounting past experiences can also result in more fluctuations as workers are less locked-in and wages in response would be more volatile.

Proposition 5. (*Time to “Unlearn” the Past with Partial Memory*) Upon change in wage regime at time $t = \tilde{t}$ from $G1$ to $G2$, there exists a time lag $(\bar{t} - \tilde{t})$ for workers to adapt from learning NE_1 to NE_2 . The transition time between games in overcoming initial propensities can be captured by the base learning rate $(\frac{1}{\eta})$ and payoff dynamics in playing $G2$:

$$(\bar{t} - \tilde{t}) \propto \frac{1}{\eta \cdot f(G2 \text{ payoff dynamics})} \quad (37)$$

Higher η , lower weight on past propensities, thus shorter time lag.

Proof. See Appendix A.2. □

In Figure 14, I show a simulation of changes in wages and choice probabilities over time when there exist a forgetting parameter (e.g. $\eta = 0.8$). As compared to perfect recall, such learning trajectory suggests that workers are less locked-in by past experiences and place greater weight on recent payoffs. They switch faster to playing the set of NE(s) supported by the wage regime they are in, thus can be perceived as more adaptive to new market conditions. While this is beneficial in stimulating switching in job applications and inducing coordination among workers when they could begin the application rounds with overcrowding at one of the firms. The pitfall of this is the volatility in choice probabilities. There could be many periods of mismatch as workers could forget previous propensities and they may not exhibit convergence to applying more to different firms even over a finite and extended period of time. (More simulation examples in Appendix B.1).

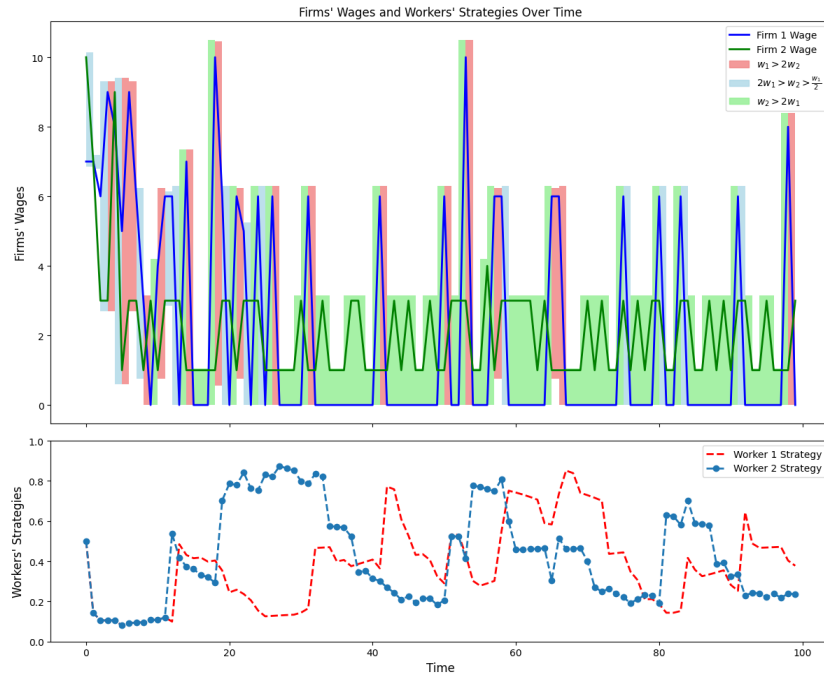


Figure 14: Changes in Wages and Workers' Choice Probabilities Over Time with $\eta = 0.8$

2.2.2 Policy Implications

Building on top of static wage environment, the dynamic wage setting portrays a slightly more realistic scenario where firms are also adaptive learners. As workers settle into more stable strategies and firms drive wages down to 0, there is a continuum of equilibria, and equilibrium selection would depend on the learning path. However, this long run behaviour is less “interesting” as workers and firms simply stop learning due to 0 reward, it is less informative of whether the eventual strategy is in fact optimal. The merit of this learning mechanism is it allows me to explore the evolution of workers’ and firms’ behaviours when wages remain positive and dynamically changing.

Since wages could change over time, workers may face 3 possible games, consisting of different sets of NE(s). Workers’ choice probabilities could be noisier as they maneuver within different wage regimes. In order to nudge workers towards a more efficient outcome, policies could focus on maintaining a wage environment with pure NEs, where workers could converge to applying to different firms eventually. Furthermore, since initial experiences could create inertia for workers to switch to learning a different set of NE(s), if workers persistently overcrowd at firm 1 for a substantial number of periods, this will influence future behaviours of workers and affect their convergence rate. This could also disproportionately benefit firm that happens to be setting higher wages at the initial application rounds. Policy makers looking to decrease probability of overcrowding could consider measures for workers to “unlearn” past experiences. While forgetting can happen naturally, it is also possible to de-emphasize past instances by re-framing each application rounds as less analogous to one another and highlighting that market conditions has changed drastically from a distant past. Imposing higher discounting could improve transition to learning a new set of NE(s) that may be more efficient, but there is a trade-off between higher adaptivity to new situations and increased volatility of choice probabilities.

3 Search with Sequential Learning

In this section, I explore the second market structure, which resembles job search on online platforms where the wages are fully revealed. Firms post wages first before workers choose the firm to apply to. Workers are assumed to adopt a logit choice model and follows best-response (BR) dynamics highlighted by Fudenberg and Levine (1998) and Hopkins (1999). They consciously study the wage environment, and respond to the perceived strategy of their opponent given their experiences. This is then feedback into firms' decision problem, where they set wages knowing how workers formulate their strategies.

The model tracks two possible sources of experiences. The first constitutes of what workers observe their opponent did over time. In directed search literature, such as Wright et al. (2021), wages often have a role of directing workers, higher wage would attract higher application rate. However, when wages are the same, workers could be equally attracted by both firms, and experiences that reflect the opponent's strategy could help determine where workers coordinate on. Even when wages differ, this influence does not vanish. Therefore, exploring the role of experiences may shine a light on equilibrium selection in presence of multiple equilibria, and also, potentially offer an explanation for why workers hardly switch in applying for different jobs despite knowing wages beforehand.

The other source could come from workers having accumulated experiences that forms a static bias that is exogenous to the current learning problem. This can be perceived as an anchoring bias based on historical events (Lieder et al. (2018)), or as social and cultural stereotypes (Langenhove and Harré (1994)) that are less shaped by short-term encounters. In Barbulescu and Bidwell (2013), they show that similarly qualified students in MBA program were found to display gender segregation in job applications, where women are less likely to apply for traditionally masculine jobs than men due to gender role stereotypes. Therefore, workers can be perceived to possess bias from long-term experiences, which affect their application choices.

Under this framework, I seek to explore how past experiences influence workers' job search behaviour, given that they observe the wages, and how this affects equilibrium selection as compared to the previous market structure.

3.1 Experience-driven Job Search with Best Response

In a similar 2×2 set-up with 2 firms and 2 workers, where workers are indexed by $i = 1, 2$, firms by $j = 1, 2$, and each firm offering one vacancy.

Timing and Set-up. Firms observe exogenous realization of productivity, $\mathbf{z} = \{z_1, z_2\}$, $\mathbf{z} \in Z$, which holds throughout all the periods. Time is discrete, $t = \{0, 1, \dots, T\}$, representing generations of workers and firms.³ In each period:

1. Firms infer workers' belief based on observed history of actions, and they set wages, $\mathbf{w} = \{w_{1t}, w_{2t}\}$, in respect to how workers are expected to react in each period.
2. Workers observe wages and choose their application strategies, given beliefs about opponent's choices and any bias they adopted from the same indexed worker of the past period.⁴
3. All parties observe the payoffs at the end of the period. Workers update their beliefs based on realized actions, and they drop out of the market, never to return.

Workers' side.

³Same as the previous set-up, I consider $T \rightarrow \infty$, which implies infinitely many generations.

⁴Same as previous set-up, worker 1 of $t = 1$ learns from worker 1 of $t = 0$, same applies for worker 2.

- **Actions.** Both workers are choosing between firm 1 and 2, $A^i = A^{-i} = \{F1, F2\}$.
- **Strategies.** Let $x = (x_1, x_2)$, $y = (y_1, y_2)$, where (x_1, x_2) and (y_1, y_2) denote the probability distribution over the two actions for worker 1 and 2 respectively. At time t , their choice probabilities are: $x_t = (x_{1t}, x_{2t})$ and $y_t = (y_{1t}, y_{2t})$, where $\sum x_t = 1$, $\sum y_t = 1$.
- **Beliefs.** Workers do not observe the exact strategies of their opponent, they form beliefs about the other worker's choice probability through realized actions. Worker 2's belief about worker 1's choice probabilities is denoted as $u_t = (u_{1t}, 1 - u_{1t})$, where $u_{1t} \in (0, 1)$; and worker 1's belief about worker 2's choice probabilities is $v_t = (v_{1t}, 1 - v_{1t})$, $v_{1t} \in (0, 1)$.
- **Bias.** Worker 1's bias is $\alpha^i = (\alpha_1^i, \alpha_2^i)$, where α_1^i is the bias towards selecting firm 1, and α_2^i is the bias towards selecting firm 2. Correspondingly, worker 2's bias is $\alpha^{-i} = (\alpha_1^{-i}, \alpha_2^{-i})$.
- **Expected Payoffs.** Given beliefs about the other worker's choice probability, the expected payoffs for worker 1 and 2 when selecting an action at time t :

$$\pi_t^i(a_t^i, v_t) = W_t(a_t^i, F1)v_{1t} + W_t(a_t^i, F2)v_{2t}, \text{ where } a_t^i \in A^i \quad (38)$$

$$\pi_t^{-i}(a_t^{-i}, u_t) = W_t^T(a_t^{-i}, F1)u_{1t} + W_t^T(a_t^{-i}, F2)u_{2t}, \text{ where } a_t^{-i} \in A^i \quad (39)$$

where the payoff matrices are

$$W_t = \begin{pmatrix} \frac{w_{1t}}{2} & w_{1t} \\ w_{2t} & \frac{w_{2t}}{2} \end{pmatrix}, W_t^T = \begin{pmatrix} \frac{w_{1t}}{2} & w_{2t} \\ w_{1t} & \frac{w_{2t}}{2} \end{pmatrix} \quad (40)$$

- **Choice Rule.** Workers are reacting to wages, their beliefs about the other worker's strategy and their own bias. β is a rationality or sensitivity parameter to the expected payoffs. For worker 1:

$$x_t = BR_x(w_{1t}, w_{2t}, v_t) = \begin{cases} x_{1t} = \frac{\exp(\alpha_1^i + \beta\pi_t^i(F1, v_t))}{\exp(\alpha_1^i + \beta\pi_t^i(F1, v_t)) + \exp(\alpha_2^i + \beta\pi_t^i(F2, v_t))} \\ x_{2t} = \frac{\exp(\alpha_2^i + \beta\pi_t^i(F2, v_t))}{\exp(\alpha_1^i + \beta\pi_t^i(F1, v_t)) + \exp(\alpha_2^i + \beta\pi_t^i(F2, v_t))} = 1 - x_{1t} \end{cases} \quad (41)$$

For worker 2:

$$y_t = BR_y(w_{1t}, w_{2t}, u_t) = \begin{cases} y_{1t} = \frac{\exp(\alpha_1^{-i} + \beta\pi_t^{-i}(F1, u_t))}{\exp(\alpha_1^{-i} + \beta\pi_t^{-i}(F1, u_t)) + \exp(\alpha_2^{-i} + \beta\pi_t^{-i}(F2, u_t))} \\ y_{2t} = \frac{\exp(\alpha_2^{-i} + \beta\pi_t^{-i}(F2, u_t))}{\exp(\alpha_1^{-i} + \beta\pi_t^{-i}(F1, u_t)) + \exp(\alpha_2^{-i} + \beta\pi_t^{-i}(F2, u_t))} = 1 - y_{1t} \end{cases} \quad (42)$$

- **Expected Motion.** Learning dynamics are captured by changes in choice probabilities:

$$\mathbb{E}(x_{t+1}) - x_t = BR_x(w_{1t}, w_{2t}, v_t) - x_t, \mathbb{E}(y_{t+1}) - y_t = BR_y(w_{1t}, w_{2t}, u_t) - y_t \quad (43)$$

- **Updating Rule.** Following Hopkins (2002), workers' beliefs about their opponent's choice probabilities are updated in each period after observing the realized actions by both workers in the previous period, and the weight attributed to initial beliefs is assumed to be 1.

$$u_{t+1} = \frac{(t+1)u_t + a_t^i}{t+2}, v_{t+1} = \frac{(t+1)v_t + a_t^{-i}}{t+2} \quad (44)$$

This can be expressed as running average of actions chosen in each period:

$$u_t = \frac{1}{t} \sum_{k=1}^t a_k^i, v_t = \frac{1}{t} \sum_{k=1}^t a_k^{-i} \quad (45)$$

In period $t+1$, payoff matrices will also be adjusted based on firms' inference about workers' behaviours. Therefore, workers' expected payoffs are updated given new set of payoff matrices, (W_{t+1}, W_{t+1}^T) , and beliefs, (v_{t+1}, u_{t+1}) .

For the workers' side, the subgame equilibrium solution(s) resembles that of quantal response equilibrium (QRE), where workers eventually form correct beliefs about opponents' strategies (McKelvey and Palfrey (1995)). Given a set of wages (w_1, w_2) , the presence of multiple equilibria depends on the sensitivity parameter, β . Higher β implies one is more sensitive to changes in expected payoffs and there is less noise in strategies, thus close to pure strategies could exist; whereas as β tends to 0, a unique mixed equilibrium emerge. Suppose workers do not possess any bias (i.e. $\alpha^i = \alpha^{-i} = 0$), then when firms are homogenous and wages equalize, (i.e. $w_1 = w_2$, $w_1 > 0$, $w_2 > 0$), Figure 25 illustrates workers' behaviours given variations in β . It shows that there can be multiple equilibria.

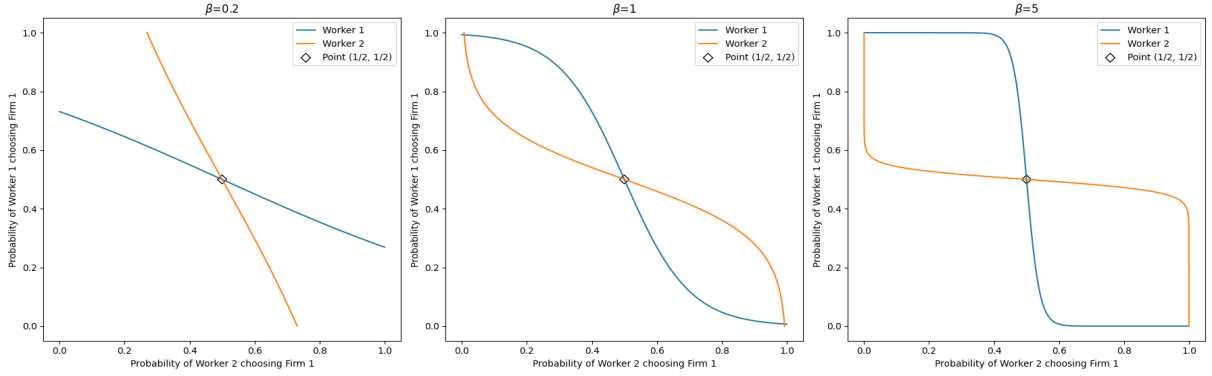


Figure 15: Workers' Response Functions to Given Set of Wages

Firms' side. Firms keeps a tally of workers' observed actions and infer workers' beliefs, they set wages in each period knowing how workers would respond. Their maximization problems:

$$\max_{w_{1t}} (1 - (1 - x_{1t})(1 - y_{1t}))(z_1 - w_{1t}) \text{ s.t. } z_1 \geq w_{1t} \geq 0 \quad (46)$$

$$\max_{w_{2t}} (1 - x_{1t}y_{1t})(z_2 - w_{2t}) \text{ s.t. } z_2 \geq w_{2t} \geq 0 \quad (47)$$

Firms' FOCs:

$$1 - (1 - x_{1t})(1 - y_{1t}) = (z_1 - w_{1t})[(1 - y_{1t})\frac{dx_{1t}}{dw_{1t}} + (1 - x_{1t})\frac{dy_{1t}}{dw_{1t}}] \quad (48)$$

$$1 - x_{1t}y_{1t} = -(z_2 - w_{2t})[y_{1t}\frac{dx_{1t}}{dw_{2t}} + x_{1t}\frac{dy_{1t}}{dw_{2t}}] \quad (49)$$

Workers' FOCs:

$$\frac{dx_{1t}}{dw_{1t}} = x_{1t}(1 - x_{1t})(\beta\frac{1}{2}v_{1t} + \beta v_{2t}) \quad (50)$$

$$\frac{dx_{1t}}{dw_{2t}} = -x_{1t}(1 - x_{1t})(\beta v_{1t} + \beta\frac{1}{2}v_{2t}) \quad (51)$$

$$\frac{dy_{1t}}{dw_{1t}} = y_{1t}(1 - y_{1t})(\beta\frac{1}{2}u_{1t} + \beta u_{2t}) \quad (52)$$

$$\frac{dy_{1t}}{dw_{2t}} = -y_{1t}(1 - y_{1t})(\beta u_{1t} + \beta\frac{1}{2}u_{2t}) \quad (53)$$

Combining them, the wage equations are:

$$w_{1t} = \max[z_1 - \frac{1 - (1 - x_{1t})(1 - y_{1t})}{(1 - y_{1t})x_{1t}(1 - x_{1t})(\beta\frac{v_{1t}}{2} + \beta(1 - v_{1t})) + (1 - x_{1t})y_{1t}(1 - y_{1t})(\beta\frac{u_{1t}}{2} + \beta(1 - u_{1t}))}, 0] \quad (54)$$

$$w_{2t} = \max[z_2 - \frac{1 - x_{1t}y_{1t}}{y_{1t}x_{1t}(1 - x_{1t})(\beta v_{1t} + \beta\frac{(1-v_{1t})}{2}) + x_{1t}y_{1t}(1 - y_{1t})(\beta u_{1t} + \beta\frac{(1-u_{1t})}{2})}, 0] \quad (55)$$

Proposition 6 (Long-run Behaviours). *Workers and firms' behaviours converge to form equilibria similar to the Quantal Response Equilibria (QRE).*

- (Workers' Beliefs.) *Workers' beliefs of opponents' choices converge to long run average of observed actions.*

$$u^* = \lim_{t \rightarrow \infty} u_t = \mathbb{E}(a^i), v^* = \lim_{t \rightarrow \infty} v_t = \mathbb{E}(a^{-i}) \quad (56)$$

- (Workers' Strategies.) *Workers' strategies converge to (x^*, y^*) for stabilized beliefs.*

$$x^* = \begin{cases} x_1 = \frac{\exp(\alpha_1^i + \beta \pi^i(F1, v^*))}{\exp(\alpha_1^i + \beta \pi^i(F1, v^*)) + \exp(\alpha_2^i + \beta \pi^i(F2, v^*))} \\ x_2 = 1 - x_1 \end{cases} \quad (57)$$

$$y^* = \begin{cases} y_1 = \frac{\exp(\alpha_1^{-i} + \beta \pi^{-i}(F1, u^*))}{\exp(\alpha_1^{-i} + \beta \pi^{-i}(F1, u^*)) + \exp(\alpha_2^{-i} + \beta \pi^{-i}(F2, u^*))} \\ y_2 = 1 - y_1 \end{cases} \quad (58)$$

where π^i and π^{-i} are expected payoffs computed using W^* , u^* and v^* .

- (Firms' Wages.) *Wages (w_1^*, w_2^*) are determined given (x^*, y^*, u^*, v^*) .*

$$w_1^* = \max[z_1 - \frac{1 - (1 - x_1^*)(1 - y_1^*)}{(1 - y_1^*)x_1^*(1 - x_1^*)(\beta \frac{v_1^*}{2} + \beta(1 - v_1^*)) + (1 - x_1^*)y_1^*(1 - y_1^*)(\beta \frac{u_1^*}{2} + \beta(1 - u_1^*))}, 0] \quad (59)$$

$$w_2^* = \max[z_2 - \frac{1 - x_1^*y_1^*}{y_1^*x_1^*(1 - x_1^*)(\beta v_1^* + \beta \frac{(1 - v_1^*)}{2}) + x_1^*y_1^*(1 - y_1^*)(\beta u_1^* + \beta \frac{(1 - u_1^*)}{2})}, 0] \quad (60)$$

and payoff matrix:

$$W^* = \begin{pmatrix} \frac{w_1^*}{2} & \frac{w_1^*}{2} \\ w_2^* & \frac{w_2^*}{2} \end{pmatrix} \quad (61)$$

- (Consistency Condition.) *In the equilibrium, $u^* = x^*$, $v^* = y^*$, beliefs matches the actual choice probabilities.*
- (Equilibrium Multiplicity.) *For given set of parameters, multiple equilibria could exist, defined by different sets of $(x_1^*, y_1^*, w_1^*, w_2^*)$.*

Proof. Based on equation (45), as $t \rightarrow \infty$, workers' beliefs of opponents' choices converges to the running average of observed actions, shown by (56). Given u^* and v^* , workers' strategies (x^*, y^*) are computed based on equations (57) and (58). Firms' equilibrium wage-setting (w_1^*, w_2^*) given workers' response is determined by equations (54) and (55).

In the equilibrium, belief consistency needs to be achieved, such that their beliefs match the choice probabilities. The system converge to equilibrium defined by $(x_1^*, y_1^*, w_1^*, w_2^*)$ as $t \rightarrow \infty$:

$$x_1^* = \frac{\exp(\alpha_1^i + \beta \pi^i(F1, y^*))}{\exp(\alpha_1^i + \beta \pi^i(F1, y^*)) + \exp(\alpha_2^i + \beta \pi^i(F2, y^*))} \quad (62)$$

$$y_1^* = \frac{\exp(\alpha_1^{-i} + \beta \pi^{-i}(F1, x^*))}{\exp(\alpha_1^{-i} + \beta \pi^{-i}(F1, x^*)) + \exp(\alpha_2^{-i} + \beta \pi^{-i}(F2, x^*))} \quad (63)$$

$$w_1^* = \max[z_1 - \frac{1 - (1 - x_1^*)(1 - y_1^*)}{(1 - y_1^*)x_1^*(1 - x_1^*)(\beta \frac{y_1^*}{2} + \beta(1 - y_1^*)) + (1 - x_1^*)y_1^*(1 - y_1^*)(\beta \frac{x_1^*}{2} + \beta(1 - x_1^*))}, 0] \quad (64)$$

$$w_2^* = \max[z_2 - \frac{1 - x_1^*y_1^*}{y_1^*x_1^*(1 - x_1^*)(\beta y_1^* + \beta \frac{(1 - y_1^*)}{2}) + x_1^*y_1^*(1 - y_1^*)(\beta x_1^* + \beta \frac{(1 - x_1^*)}{2})}, 0] \quad (65)$$

These equations collectively define the equilibrium. Multiple solutions could exist. \square

In the following Algorithm 1, I show firms' wage-setting behaviours using a grid search approach, where firms conduct a coarse search and then a refined search to nail down the wages in each period based on potential worker reactions.

Algorithm 1 Firms' Wage-setting using Grid Search Approach

```

1: Initialize for  $t = 0$ , set  $\alpha_1^i, \alpha_1^{-i}, u_0^i, v_0^i, x_0, y_0$ ; Compute  $w_{10}, w_{20}$ .
2: for one session do
3:   Loop the following
4:   for 100000 time periods do
5:     Loop for each time period
6:     for all firms do
7:       Conduct coarse search, followed by more refined search
8:       for coarse search do
9:         Set two arrays of wages bounded by  $[0, z_j]$ , divide the range into 10 evenly
           spaced values, such that there can be finite pairs of  $(w_{1t}^{\text{Coarse}}, w_{2t}^{\text{Coarse}})$ .
10:        Compute workers' reaction based on equations (41) and (42) to obtain
            $(x_{1t}^{\text{Potential}}, y_{1t}^{\text{Potential}})$  to each pair of  $(w_{1t}^{\text{Coarse}}, w_{2t}^{\text{Coarse}})$ .
11:        Based on workers' potential application rates, compute firms' payoffs using
           equations (46) and (47).
12:        Find the wage pair that lead to highest profit.
13:      end for coarse search
14:      for refined search do
15:        Take the pair of wages previously identified,  $(w_{1t}^{\text{Coarse}}, w_{2t}^{\text{Coarse}})$ , and create a
           finer search window (i.e.  $\pm 0.5$ ), dividing this range into 20 equal units.
16:        Search all combination of wages in this range to find the pair that maximizes
           individual firm's profit,  $(w_{1t}^{\text{Refined}}, w_{2t}^{\text{Refined}})$ .
17:      end for refined search
18:      Set the wages  $(w_{1t}, w_{2t}) = (w_{1t}^{\text{Refined}}, w_{2t}^{\text{Refined}})$ .
19:    end for firms
20:    for all workers do
21:      Observe  $(w_{1t}, w_{2t})$ , compute  $x_t$  and  $y_t$  given  $u_t$  and  $v_t$  using equations (41) and
           (42).
22:      Generate a choice of action from  $(F1, F2)$  for each worker based on  $x_t$  and  $y_t$ .
23:    end for workers
24:    for reward generation and updating do
25:      Given realized workers' choices,  $(a_t^i, a_t^{-i})$ , and firms' wages  $(w_{1t}, w_{2t})$ , compute
           the rewards for all agents.
26:      Workers' beliefs about each other  $(u_{t+1}, v_{t+1})$  are updated based on equation
           (44) for use in the following period.
27:    end for one period
28:  end for all periods
29: end for all sessions
30: return results  $(x_{1t}, x_{2t}, y_{1t}, y_{2t}, u_{1t}, u_{2t}, v_{1t}, v_{2t}, w_{1t}, w_{2t})$  for all periods.

```

In Figure 16, I show the convergence pattern of wages (top panel), choice probabilities (middle panel) and beliefs (bottom panel) given different β s. When $\beta = 0.2$, there is a unique equilibrium, and there is convergence to it. In presence of multiple equilibria, it is shown for $\beta = 1.0$, workers converge to a strategy that is close to random search; and for $\beta = 5.0$, workers' converge to

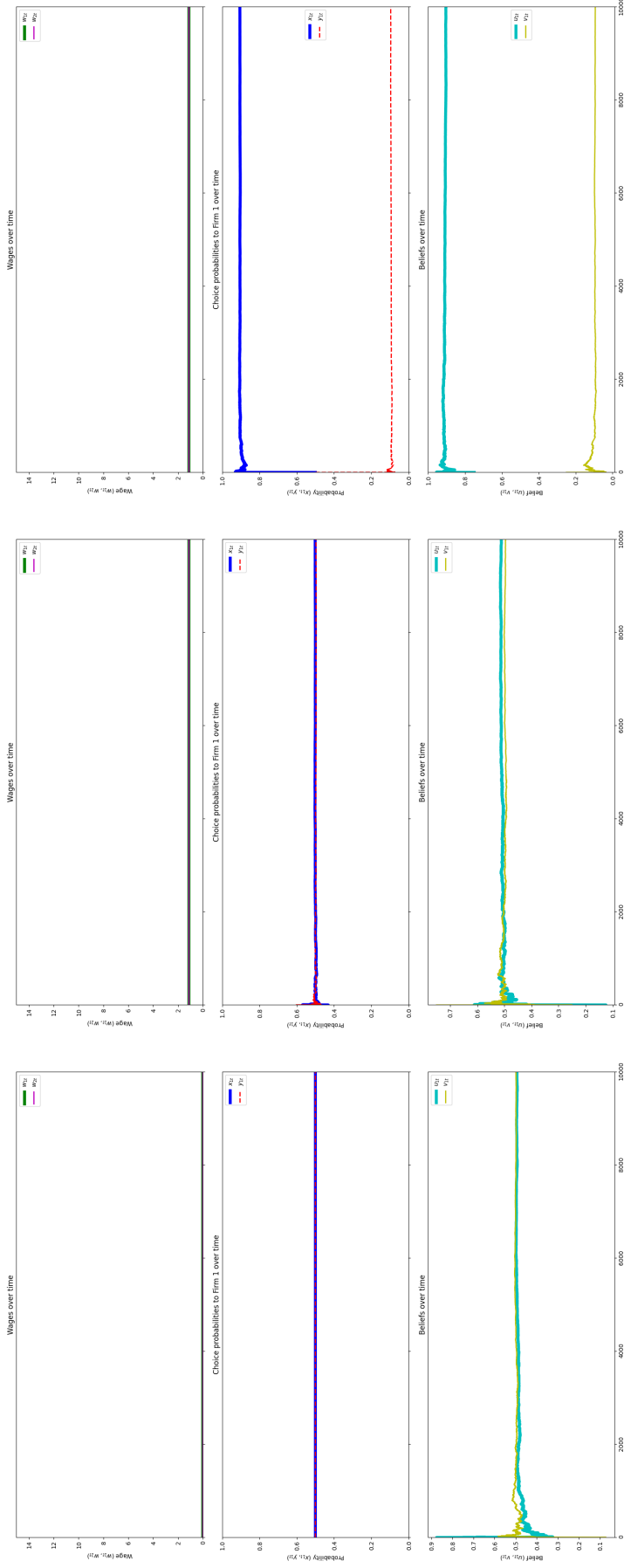


Figure shows wages, choice probabilities and belief evolution for different $\beta = \{0.2, 1.0, 5.0\}$, given $z_1 = z_2 = 10$, $t = 10000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0.5$, $v_{10} = 0.5$, no smoothing.

Figure 16: Changes in Wages and Workers' Choice Probabilities Over Time (Algorithm 1)

almost pure coordination on applying with high probability to different firms. Since this grid search wage-setting method emphasizes on maximum profits for the firms, if wages differ across equilibria, the firms would converge towards an equilibrium with lower wages. The simulation results draws some resemblance to the theoretical findings in Lu (2024), where at relatively high sensitivity to expected payoffs, wages for asymmetric equilibria tend to correspond to lower wages as compared to symmetric equilibrium; and vice versa for slightly lower sensitivity to expected payoffs. As a result, when firms maximize their profits by setting lower wages, workers would converge to the asymmetric strategies for $\beta = 5.0$, and to symmetric strategies for $\beta = 1.0$, given that multiple equilibria exist. Whilst simulation results in Figure 16c show convergence to a single asymmetric equilibrium under this approach, it is necessary to note that the two asymmetric equilibria are equally attractive, and which one is selected would be contingent on the learning path.

An important distinction between the equilibrium achieved as compared to traditional QRE is that workers and firms are not forward-looking. Workers are myopic in a sense that they only update their beliefs based on past observations of opponents' actions. They do not strategize given how their current actions might influence the other worker's behaviour. For the firms, they also have limited foresight. While they take into account workers' reaction to the wages posted in the current period and optimize based on potential response, they do not account for how wages might affect beliefs. Therefore, this myopic learning dynamics may converge only to a subset of QRE that can be identified in forward-looking scenarios.

To track the different parts of the system, I first fix workers' beliefs to investigate how would choice probabilities and wages behave; I then fix wages to explore how might beliefs and choice probabilities evolve. In Figure 17, I show that when beliefs are fixed, workers will converge to a unique set of choice probabilities. In presence of multiple equilibria, the point at which beliefs stabilize is important in determining if one converge to the asymmetric or the symmetric equilibrium. If workers hold beliefs that the opponent applies more to a different firm than them, then it is more likely for workers to converge to applying asymmetrically. In view of workers' behaviours stabilize, wages also stabilize. In Figure 18, I fix either both of the wages or one of the wages. In both cases, workers' beliefs becomes asymmetric and they converge to applying more to different firms.

Based on Proposition 6 and simulations, when β is low and only symmetric equilibrium exists, there is convergence to the unique equilibrium; when β is high and multiple equilibria exist, there will be convergence to the asymmetric equilibria. This is more efficient than the symmetric case due to higher chances of one-to-one matching. However, which of the two asymmetric equilibria the workers coordinate on could depend heavily on beliefs.

Apart from grid search approach, I have also attempted other wage-setting algorithms, shown in Appendix B.2. The grid search approach could implement a cleaner selection of one of the asymmetric equilibria, but would require firms to have perfect information of the profit landscape for any wage combinations, as well as the computational power to evaluate the expected payoffs. However, wages tend to be low as firms search through all combinations of wages to maximize profits, which implies monopsony power can be high.

Proposition 7 (Equilibrium Stability). *Let $(x_1^*, y_1^*, w_1^*, w_2^*)$ be an equilibrium point in the firm-worker matching problem. The dynamic system is defined by:*

$$\dot{u}_t = x_t(v_t; w_{1t}, w_{2t}) - u_t \quad (66)$$

$$\dot{v}_t = y_t(u_t; w_{1t}, w_{2t}) - v_t \quad (67)$$

where x_t, y_t, w_{1t}, w_{2t} depends on u_t and v_t . In equilibrium, $u^* = x^*, v^* = y^*$. The Jacobian of

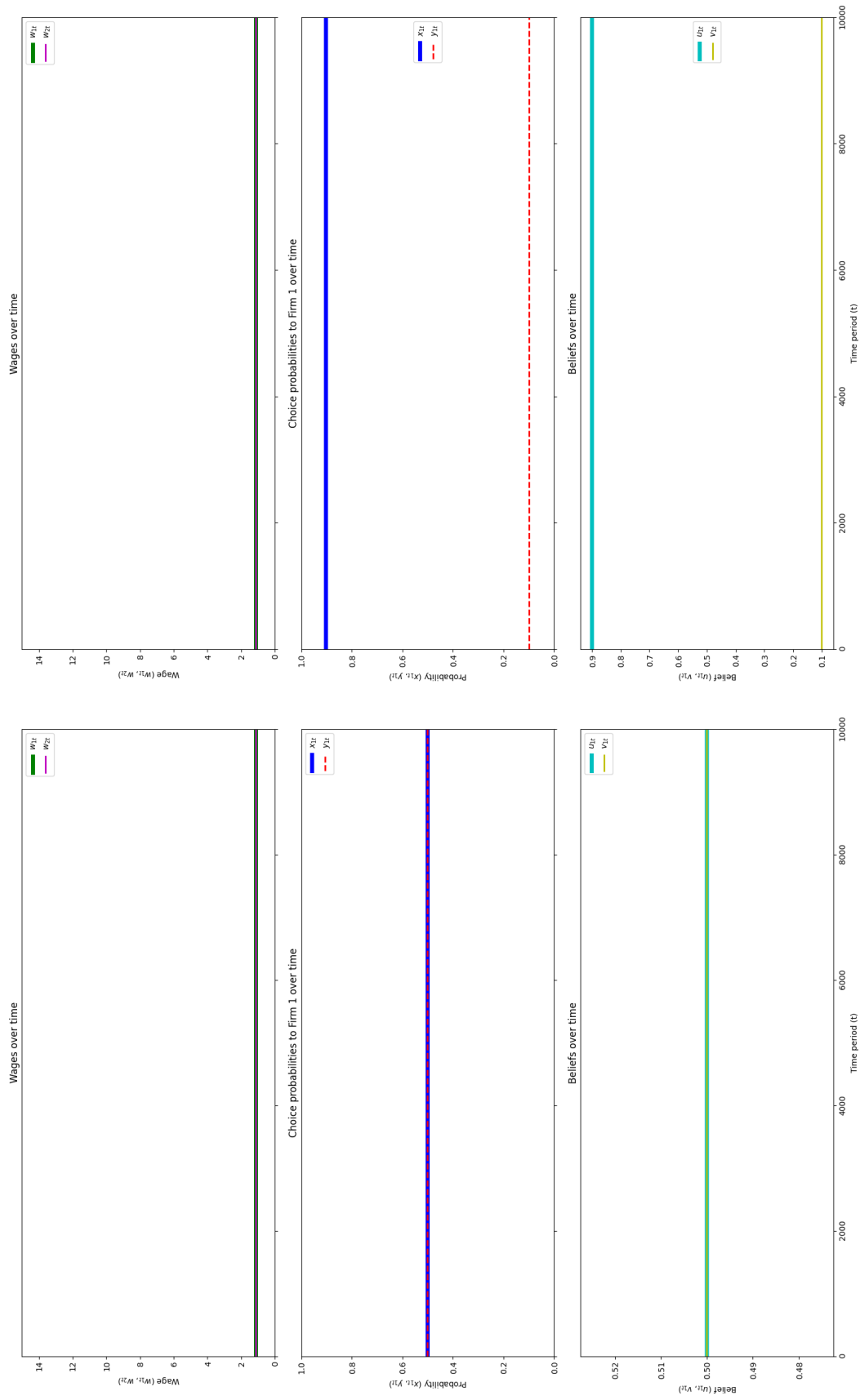
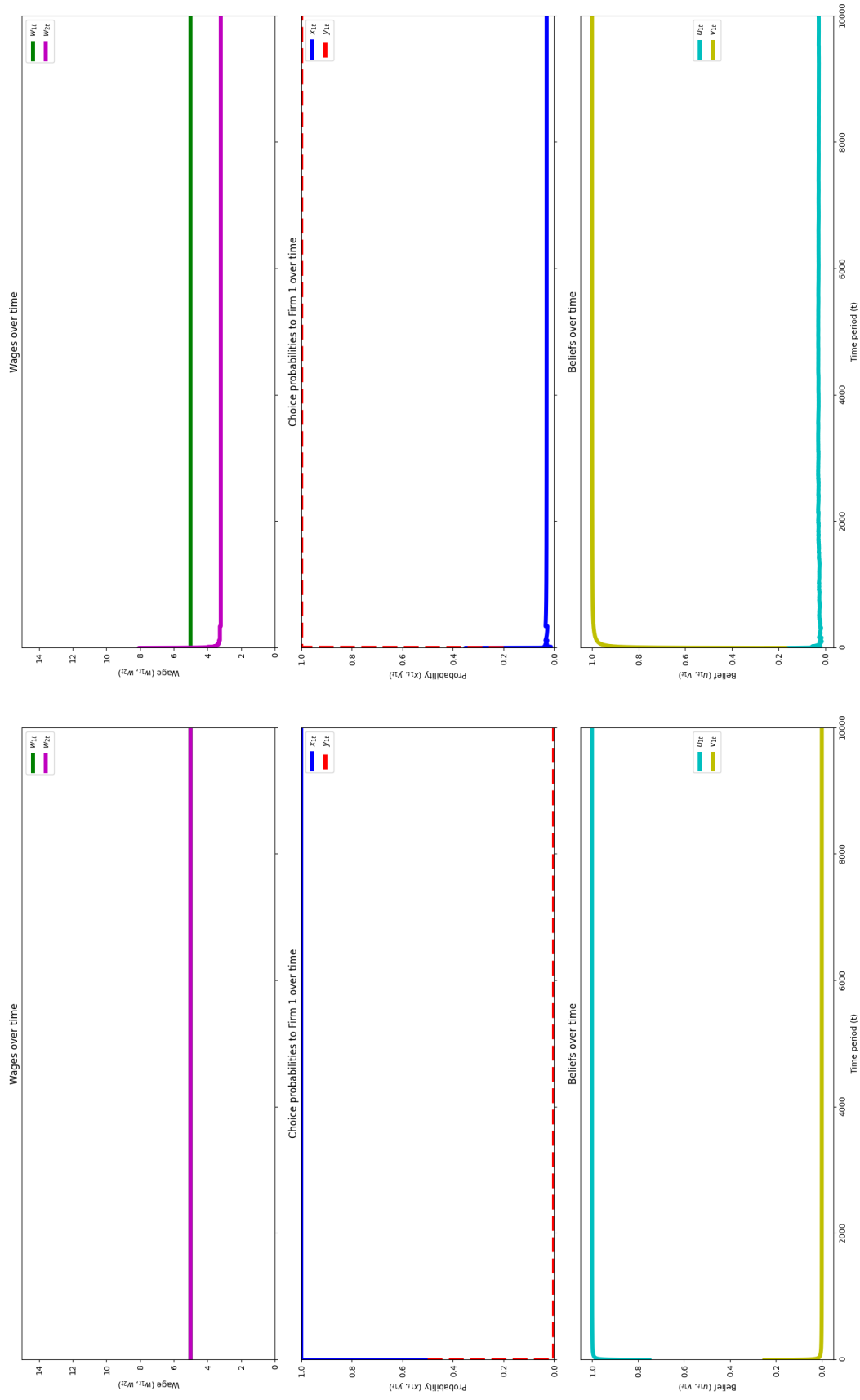


Figure shows wages, choice probabilities, fixing beliefs at $(u_1, v_1) = (0.5, 0.5)$ and $(0.9, 0.1)$, given $\beta = 5.0$, $z_1 = z_2 = 10$, $t = 10000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0$, and no smoothing.

Figure 17: Changes in Wages, Workers' Choice Probabilities Over Time for Fixed Beliefs



(a) $w_1 = w_2 = 5$

(b) $w_1 = 5$

Figure shows wages, choice probabilities and belief evolution, given $\beta = 5.0$, $t = 10000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0$, $u_{10} = 0.5$, $v_{10} = 0.5$, no smoothing.

Figure 18: Changes in Workers' Choice Probabilities, Beliefs Over Time for Fixed Wages

the system evaluated at the equilibrium points is given by:

$$J = J_{base} + W_1 + W_2 \quad (68)$$

where the full form is:

$$J = \begin{pmatrix} 0 & \frac{\partial x_t}{\partial v_t} \\ \frac{\partial y_t}{\partial u_t} & 0 \end{pmatrix} - I + \begin{pmatrix} \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} & \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} \\ \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} & \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} \end{pmatrix} + \begin{pmatrix} \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t} & \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t} \\ \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t} & \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t} \end{pmatrix} \quad (69)$$

J_{base} captures the direct effect of u_t and v_t on x_t and y_t ; and W_1, W_2 account for the indirect effect of u_t and v_t on x_t and y_t through w_{1t} and w_{2t} . In the equilibrium,

- For $W_1 + W_2 \approx 0$, asymmetric equilibria, if exist, are locally asymptotically stable; symmetric equilibrium is a saddle point if β is large and stable if β is small.
- If $W_1 + W_2$ is not small relative to J_{base} , the indirect effect may impact the signs of the eigenvalues and potentially destabilize the equilibrium.

Corollary 3.0.1 (Stability with Approximately Constant Wages.). *If w_1, w_2 are approximately constant, then $W_1 + W_2 \approx 0$, the eigenvalues of J_{base} :*

$$\lambda = -1 \pm \sqrt{\frac{\partial x_t}{\partial v_t} \frac{\partial y_t}{\partial u_t}} \quad (70)$$

If $\frac{\partial x_t}{\partial v_t} \frac{\partial y_t}{\partial u_t} < 1$, which holds in the limit of $x_1, y_1 \rightarrow (0, 1)$ or $(1, 0)$, then all eigenvalues have negative real parts, and the system is locally asymptotically stable.

Proof. In order to analyze stability of the equilibrium outcome(s). I restate the following equations for clearer view.

Beliefs:

$$u_{t+1} = \frac{(t+1)u_t + a_t^i}{t+2}, v_{t+1} = \frac{(t+1)v_t + a_t^{-i}}{t+2} \quad (71)$$

Workers' choice probabilities:

$$x_{1t} = \frac{1}{1 + \exp(-[(\alpha_1^i - \alpha_2^i) + \beta(w_{1t} - \frac{w_{2t}}{2} - (\frac{w_{1t}}{2} + \frac{w_{2t}}{2})v_{1t})])} \quad (72)$$

$$y_{1t} = \frac{1}{1 + \exp(-[(\alpha_1^{-i} - \alpha_2^{-i}) + \beta(w_{1t} - \frac{w_{2t}}{2} - (\frac{w_{1t}}{2} + \frac{w_{2t}}{2})u_{1t})])} \quad (73)$$

Wages:

$$w_{1t} = \max[z_1 - \frac{1 - (1 - x_{1t})(1 - y_{1t})}{(1 - y_{1t})x_{1t}(1 - x_{1t})(\beta\frac{v_{1t}}{2} + \beta(1 - v_{1t})) + (1 - x_{1t})y_{1t}(1 - y_{1t})(\beta\frac{u_{1t}}{2} + \beta(1 - u_{1t}))}, 0] \quad (74)$$

$$w_{2t} = \max[z_2 - \frac{1 - x_{1t}y_{1t}}{y_{1t}x_{1t}(1 - x_{1t})(\beta v_{1t} + \beta\frac{(1-v_{1t})}{2}) + x_{1t}y_{1t}(1 - y_{1t})(\beta u_{1t} + \beta\frac{(1-u_{1t})}{2})}, 0] \quad (75)$$

For the beliefs, the stochastic system is given by:

$$\dot{u}_t = a_t^i - u_t, \dot{v}_t = a_t^{-i} - v_t \quad (76)$$

where $a_t^i \sim \text{Bernoulli}(x_t)$, such that $a_t^i = 1$ with probability x_{1t} , and $a_t^i = 0$ with probability $1 - x_{1t}$; and $a_t^{-i} \sim \text{Bernoulli}(y_t)$.

For stability analysis, the deterministic approximation of the dynamics:

$$\dot{u}_t = \mathbb{E}(a_t^i) - u_t = x_t(v_t; w_{1t}, w_{2t}) - u_t \quad (77)$$

$$\dot{v}_t = \mathbb{E}(a_t^{-i}) - v_t = y_t(u_t; w_{1t}, w_{2t}) - v_t \quad (78)$$

where wages boil down to $w_{1t} = w_{1t}(u_t, v_t)$, $w_{2t} = w_{2t}(u_t, v_t)$. Beliefs are the only dynamical values, the choice probabilities and wages are instantaneous functions of u_t and v_t .

Find the Jacobian:

$$J = \begin{pmatrix} \frac{\partial \dot{u}_t}{\partial u_t} & \frac{\partial \dot{u}_t}{\partial v_t} \\ \frac{\partial \dot{v}_t}{\partial u_t} & \frac{\partial \dot{v}_t}{\partial v_t} \end{pmatrix} \quad (79)$$

$$\frac{\partial \dot{u}_t}{\partial u_t} = \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} + \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t} - 1 \quad (80)$$

$$\frac{\partial \dot{u}_t}{\partial v_t} = \frac{\partial x_t}{\partial v_t} + \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} + \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t} \quad (81)$$

$$\frac{\partial \dot{v}_t}{\partial u_t} = \frac{\partial y_t}{\partial u_t} + \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} + \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t} \quad (82)$$

$$\frac{\partial \dot{v}_t}{\partial v_t} = \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} + \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t} - 1 \quad (83)$$

$$J = \underbrace{\begin{pmatrix} 0 & \frac{\partial x_t}{\partial v_t} \\ \frac{\partial y_t}{\partial u_t} & 0 \end{pmatrix}}_{J_{\text{base}}} - I + \underbrace{\begin{pmatrix} \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} & \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} \\ \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} & \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} \end{pmatrix}}_{W_1} + \underbrace{\begin{pmatrix} \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t} & \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t} \\ \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t} & \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t} \end{pmatrix}}_{W_2} \quad (84)$$

If $W_1 + W_2 \approx 0$, the eigenvalues for J_{base} :

$$\lambda = -1 \pm \sqrt{\beta^2 x_{1t}(1-x_{1t})y_{1t}(1-y_{1t})\left(\frac{w_{1t}}{2} + \frac{w_{2t}}{2}\right)^2} \quad (85)$$

When β is high, multiple equilibria exist. For asymmetric equilibria (i.e. $(x_{1t}, y_{1t}) \rightarrow (0, 1)$ or $(1, 0)$ as $t \rightarrow \infty$), then $\beta(\frac{w_{1t}}{2} + \frac{w_{2t}}{2})\sqrt{x_{1t}(1-x_{1t})y_{1t}(1-y_{1t})} < 1$, the equilibria are locally asymptotically stable. For symmetric equilibrium (i.e. $x_{1t} = y_{1t}$, $x_{1t}, y_{1t} \in (0, 1)$), $\beta(\frac{w_{1t}}{2} + \frac{w_{2t}}{2})\sqrt{x_{1t}(1-x_{1t})y_{1t}(1-y_{1t})} > 1$, the equilibrium is unstable.

When β is low, only symmetric equilibrium exist. Since $\beta(\frac{w_{1t}}{2} + \frac{w_{2t}}{2})\sqrt{x_{1t}(1-x_{1t})y_{1t}(1-y_{1t})} < 1$, symmetric equilibrium is locally asymptotically stable.

The above holds true only if $W_1 + W_2$ are small relative to J_{base} .

For Jacobian matrix (69),

$$\det \begin{pmatrix} -1 - \lambda + \underbrace{\frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} + \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t}}_a & \underbrace{\frac{\partial x_t}{\partial v_t} + \frac{\partial x_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} + \frac{\partial x_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t}}_b \\ \underbrace{\frac{\partial y_t}{\partial u_t} + \frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial u_t} + \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial u_t}}_d & -1 - \lambda + \underbrace{\frac{\partial y_t}{\partial w_{1t}} \frac{\partial w_{1t}}{\partial v_t} + \frac{\partial y_t}{\partial w_{2t}} \frac{\partial w_{2t}}{\partial v_t}}_b \end{pmatrix} = 0 \quad (86)$$

$$\lambda = -1 - \frac{(a+b) \pm \sqrt{(a+b)^2 + 4(cd-ab)}}{2} \quad (87)$$

Breaking down:

$$\frac{\partial x_{1t}}{\partial w_{1t}} = \beta x_{1t}(1-x_{1t})(1 - \frac{v_{1t}}{2}) \quad (88)$$

$$\frac{\partial x_{1t}}{\partial w_{2t}} = -\beta x_{1t}(1-x_{1t})\left(\frac{1}{2} + \frac{v_{1t}}{2}\right) \quad (89)$$

$$\frac{\partial y_{1t}}{\partial w_{1t}} = \beta y_{1t}(1-y_{1t})\left(1 - \frac{u_{1t}}{2}\right) \quad (90)$$

$$\frac{\partial y_{1t}}{\partial w_{2t}} = -\beta y_{1t}(1-y_{1t})\left(\frac{1}{2} + \frac{u_{1t}}{2}\right) \quad (91)$$

$$\frac{\partial x_{1t}}{\partial v_{1t}} = -\beta x_{1t}(1-x_{1t})\left(\frac{w_{1t}}{2} + \frac{w_{2t}}{2}\right) \quad (92)$$

$$\frac{\partial y_{1t}}{\partial u_{1t}} = -\beta y_{1t}(1-y_{1t})\left(\frac{w_{1t}}{2} + \frac{w_{2t}}{2}\right) \quad (93)$$

$$\frac{\partial w_{1t}}{\partial u_{1t}} = -\frac{1}{2}\beta(1-x_{1t})y_{1t}(1-y_{1t})\frac{1-(1-x_{1t})(1-y_{1t})}{[(1-y_{1t})x_{1t}(1-x_{1t})\beta(1-\frac{v_{1t}}{2}) + (1-x_{1t})y_{1t}(1-y_{1t})\beta(1-\frac{u_{1t}}{2})]^2} \quad (94)$$

$$\frac{\partial w_{1t}}{\partial v_{1t}} = -\frac{1}{2}\beta(1-y_{1t})x_{1t}(1-x_{1t})\frac{1-(1-x_{1t})(1-y_{1t})}{[(1-y_{1t})x_{1t}(1-x_{1t})\beta(1-\frac{v_{1t}}{2}) + (1-x_{1t})y_{1t}(1-y_{1t})\beta(1-\frac{u_{1t}}{2})]^2} \quad (95)$$

$$\frac{\partial w_{2t}}{\partial u_{1t}} = \frac{1}{2}\beta x_{1t}y_{1t}(1-y_{1t})\frac{1-x_{1t}y_{1t}}{[y_{1t}x_{1t}(1-x_{1t})\beta(\frac{1+v_{1t}}{2}) + x_{1t}y_{1t}(1-y_{1t})\beta(\frac{1+u_{1t}}{2})]^2} \quad (96)$$

$$\frac{\partial w_{2t}}{\partial v_{1t}} = \frac{1}{2}\beta y_{1t}x_{1t}(1-x_{1t})\frac{1-x_{1t}y_{1t}}{[y_{1t}x_{1t}(1-x_{1t})\beta(\frac{1+v_{1t}}{2}) + x_{1t}y_{1t}(1-y_{1t})\beta(\frac{1+u_{1t}}{2})]^2} \quad (97)$$

By equation (87), as long as $\frac{(a+b) \pm \sqrt{(a+b)^2 + 4(cd-ab)}}{2} < 1$, the equilibrium will be locally asymptotically stable. Magnitude of β is essential in determining the existence of multiple equilibria, as well as stability for symmetric equilibrium in particular. For $W_1 + W_2$ further away from 0, the indirect effect of beliefs on wages becomes more important for equilibrium stability.

For the wages, as $u^* = x^*$, $v^* = y^*$,

$$\frac{\partial w_1^*}{\partial u^*} = -\frac{1}{2}y^*\frac{1-(1-x^*)(1-y^*)}{\beta(1-x^*)(1-y^*)[x^*(1-\frac{y^*}{2}) + y^*(1-\frac{x^*}{2})]^2} \quad (98)$$

$$\frac{\partial w_1^*}{\partial v^*} = -\frac{1}{2}x^*\frac{1-(1-x^*)(1-y^*)}{\beta(1-y^*)(1-x^*)[x^*(1-\frac{y^*}{2}) + y^*(1-\frac{x^*}{2})]^2} \quad (99)$$

$$\frac{\partial w_2^*}{\partial u^*} = \frac{1}{2}(1-y^*)\frac{1-x^*y^*}{\beta x^*y^*[(1-x^*)\frac{(1+y^*)}{2} + (1-y^*)\frac{(1+x^*)}{2}]^2} \quad (100)$$

$$\frac{\partial w_2^*}{\partial v^*} = \frac{1}{2}(1-x^*)\frac{1-x^*y^*}{\beta y^*x^*[(1-x^*)\frac{(1+y^*)}{2} + (1-y^*)\frac{(1+x^*)}{2}]^2} \quad (101)$$

For asymmetric equilibria, all the partial derivatives tend to 0, implying the impact of beliefs on wages is negligible. For symmetric equilibrium, some partial derivatives are negative, and some are positive, which could imply greater stabilization or destabilization as indirect effect of beliefs on wages would be non-zero. Given the sign of equations (88)-(97), stronger beliefs (i.e. increase u^* and v^*) may introduce complex eigenvalues due to non-zero imaginary parts. However, they could also contribute to more negative real part, indicating greater local asymptotic stability. \square

3.2 Anchoring Bias from Long-term Experiences

The next step is to consider non-zero α^i and α^{-i} . Rewriting based on equations (41) and (42):

$$x_{1t} = \frac{1}{1 + \exp(-[(\alpha_1^i - \alpha_2^i) + \beta(\pi_t^i(F1, v_t) - \pi_t^i(F2, v_t))])} \quad (102)$$

$$y_{1t} = \frac{1}{1 + \exp(-[(\alpha_1^{-i} - \alpha_2^{-i}) + \beta(\pi_t^{-i}(F1, u_t) - \pi_t^{-i}(F2, u_t))])} \quad (103)$$

These imply that relative bias matters more than the absolute bias, and introducing some bias towards one firm over the other could affect choice probabilities. (see Appendix B.3)

Proposition 8 (Impact of Long-term Experiences on Equilibrium Multiplicity). *The magnitude of relative bias towards a firm, $\Delta\alpha^i = \alpha_1^i - \alpha_2^i$, as compared to the expected payoff difference between the two options, $\Delta\pi^i = \pi^i(F1) - \pi^i(F2)$, could affect equilibrium multiplicity and selection.*

- *Strong Experience Bias ($|\Delta\alpha^i| > \beta|\Delta\pi^i|$): Convergence to unique equilibrium.*
- *Weak Experience Bias ($|\Delta\alpha^i| < \beta|\Delta\pi^i|$): Multiple equilibria exist, convergence to equilibrium that resembles quantal response equilibrium (QRE), but selection depends on belief updating.*

Proof. From equation (102) and (103), let $\Delta\pi^i = \pi_t^i(F1, v_t) - \pi_t^i(F2, v_t)$, $\Delta\pi^{-i} = \pi_t^{-i}(F1, v_t) - \pi_t^{-i}(F2, v_t)$; $\alpha^i = \alpha_1^i - \alpha_2^i$, $\alpha^{-i} = \alpha_1^{-i} - \alpha_2^{-i}$. If $|\Delta\alpha^i| > \beta|\Delta\pi^i|$, relative bias from experience would have a large influence on x_{1t} and y_{1t} , and beliefs about opponents would have less impact. Assuming $\beta|\Delta\pi^i| \rightarrow 0$,

$$x_{1t} \approx \frac{1}{1 + \exp(-(\alpha_1^i - \alpha_2^i))}, y_{1t} \approx \frac{1}{1 + \exp(-(\alpha_1^{-i} - \alpha_2^{-i}))} \quad (104)$$

x_{1t} and y_{1t} converge to constant values, similarly for w_{1t} and w_{2t} , following equations (64) and (65).

If $|\Delta\alpha^i| < \beta|\Delta\pi^i|$, relative expected payoffs would have a larger influence on x_{1t} and y_{1t} than relative bias. Assuming $|\Delta\alpha^i| \rightarrow 0$,

$$x_{1t} \approx \frac{1}{1 + \exp(-[\beta(\pi_t^i(F1, v_t) - \pi_t^i(F2, v_t))])}, y_{1t} \approx \frac{1}{1 + \exp(-[\beta(\pi_t^{-i}(F1, u_t) - \pi_t^{-i}(F2, u_t))])} \quad (105)$$

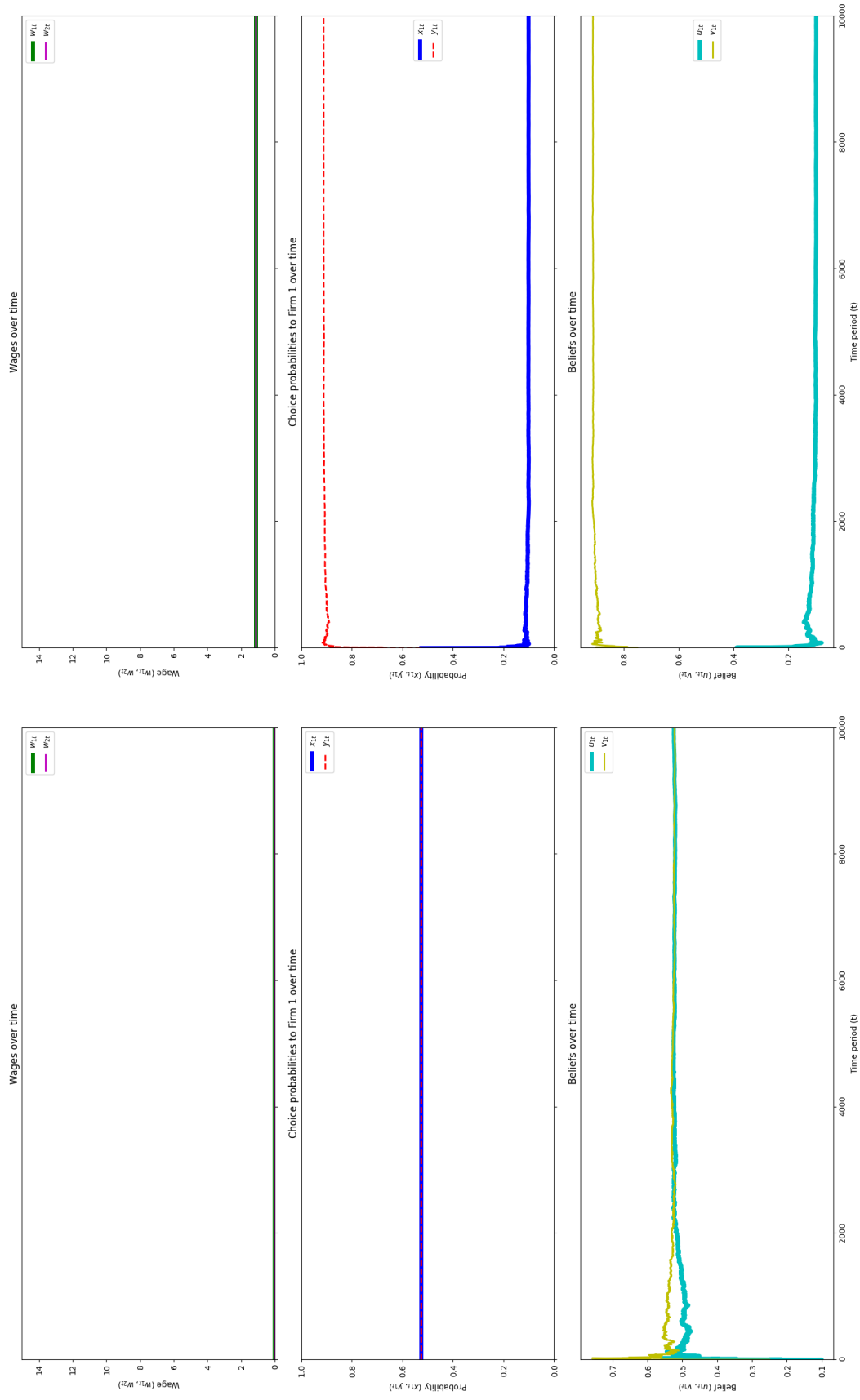
In the equilibrium where beliefs converge, the system will move towards QRE, and the selection depends on u^*, v^* . \square

Figure 19 shows simulation examples when workers have experience bias towards firm 1 relative to firm 2. There is clear convergence to one set of workers' strategies for both β s. At $\beta = 0.2$, experience bias is strong, workers' strategies are heavily influenced by it, therefore, they both apply with higher probability to firm 1. At $\beta = 5.0$, experience bias is weak, and there are multiple equilibria. Workers are more sensitive to expected payoffs. Convergence to one of the asymmetric equilibria is observed as beliefs stabilize.

From equations (72) and (73),

$$\frac{dx_{1t}}{d\Delta\alpha^i} = x_{1t}(1 - x_{1t}) > 0, \frac{dy_{1t}}{d\Delta\alpha^{-i}} = y_{1t}(1 - y_{1t}) > 0 \quad (106)$$

$\Delta\alpha^i$ and $\Delta\alpha^{-i}$ affect x_{1t} and y_{1t} positively, implying stronger bias towards firm 1 would increase application rate towards firm 1, and they would also influence wages indirectly through x_{1t} , y_{1t} .



(a) $\beta = 0.2$

(b) $\beta = 5.0$

Figure shows wages, choice probabilities and belief evolution, given $z_1 = z_2 = 10$, $t = 10000$, $\alpha_1^i = 0.1$, $\alpha_1^{-i} = 0.1$, $u_{10} = 0.5$, $v_{10} = 0.5$, no smoothing.

Figure 19: Changes in Wages and Workers' Choice Probabilities with Bias (Algorithm 1)

Corollary 3.0.2 (Stability Impact from Experience Bias under Approximately Constant Wages). *Assume $W_1 + W_2 \approx 0$, the Jacobian reduces to J_{base} , stability does not change, but stability rate is affected.*

- For asymmetric equilibria (i.e. $x_{1t}, y_{1t} \rightarrow (0, 1)$ or $(1, 0)$), impact of changes in relative experiences on eigenvalues is negligible, $\Delta\alpha^i$.
- For symmetric equilibrium (i.e. $x_{1t} = y_{1t}$), eigenvalues are affected by $\Delta\alpha^i$, there can be both stabilizing and destabilizing effect.

Proof. Based on equation (85):

$$\frac{d\lambda}{d\Delta\alpha^i} = \pm \frac{1}{2} (\beta^2 x_{1t}(1-x_{1t})y_{1t}(1-y_{1t}) (\frac{w_{1t}}{2} + \frac{w_{2t}}{2})^2)^{-\frac{1}{2}} \beta^2 (\frac{w_{1t}}{2} + \frac{w_{2t}}{2})^2 [(1-2x_{1t}) \frac{dx_{1t}}{d\Delta\alpha^i} y_{1t}(1-y_{1t}) + x_{1t}(1-x_{1t})(1-2y_{1t}) \frac{dy_{1t}}{d\Delta\alpha^i}] \quad (107)$$

Given $\frac{dx_{1t}}{d\Delta\alpha^i} > 0$, $\frac{dy_{1t}}{d\Delta\alpha^i} > 0$, for asymmetric equilibria, where $x_{1t}, y_{1t} \rightarrow (0, 1)$ or $(1, 0)$, $\frac{d\lambda}{d\Delta\alpha^i} \rightarrow 0$, the impact is negligible.

For symmetric equilibria, when β is high, the equilibrium remains unstable, but if $x_{1t} = y_{1t} > 0.5$,

$$\lambda_1 = -ve : \frac{d\lambda_1}{d\Delta\alpha^i} > 0, \lambda_2 = +ve : \frac{d\lambda_2}{d\Delta\alpha^i} < 0 \quad (108)$$

Increasing α^i could make λ_1 less negative, and λ_2 less positive. While it is still a saddle point, the trajectories may stick around the saddle for longer duration.

If $x_{1t} = y_{1t} < 0.5$,

$$\lambda_1 = -ve : \frac{d\lambda_1}{d\Delta\alpha^i} < 0, \lambda_2 = +ve : \frac{d\lambda_2}{d\Delta\alpha^i} > 0 \quad (109)$$

Increasing α^i makes λ_1 more negative, and λ_2 more positive, the stable direction is more stable, whereas the unstable direction is more unstable.

As for low β , the signs do not change so the symmetric equilibrium remain stable, but the analysis of $\frac{d\lambda}{d\Delta\alpha^i}$ applies similarly. \square

Result Summary. In this section, I explore the market structure where workers observe the wages before formulating their application strategies. Workers' and firms' behaviour converge to equilibria similar to the QRE in the long run, and there can be multiple equilibria depending on workers' sensitivity to expected payoffs (β). Using Algorithm 1, I show that as beliefs stabilize, there is clear convergence to a single equilibrium. In presence of multiple equilibria, workers and firms would converge to the symmetric equilibria given moderate β , and converge to the asymmetric equilibria when β is sufficiently high. The asymmetric equilibria, if exist, are locally asymptotically stable, while the symmetric equilibrium is only locally asymptotically stable when β is low.

The presence of exogenous experience bias could affect equilibrium multiplicity, depending on the weighing between sensitivity to expected payoffs and extent of bias. When experience bias is strong, workers would rely on their bias to inform application strategies. When experience bias is weak, multiple equilibria could emerge and previous analysis about convergence apply similarly. However, such bias does not affect equilibrium stability.

3.3 Policy Implications

In this market structure, equilibrium selection could be largely dependent on workers' sensitivity to expected payoffs, as well as their beliefs about their opponent's choice probabilities. When workers' sensitivity to expected payoff is high, workers apply with higher probabilities to different firms, contributing to more efficient outcome. When sensitivity is moderate, however, there could be convergence to the symmetric equilibrium, which is less efficient due to lower chances of one-to-one matching. As a result, measures to improve workers' sensitivity, such as through information provision, may be useful to bring attention to changes in expected payoffs, thereby inducing convergence towards the asymmetric equilibria.

When workers possess exogenous bias from long-term experiences, the extent of it against expected payoffs help to determine the existence of multiple equilibria. If workers' possess strong bias towards different firm from one another, and overriding the expected payoff difference between selecting each firm, they are directed naturally and apply with higher probability to different firms. Potential measures to implement this could be early exposure programmes, such as internships, that could prime workers to hold different bias based on early experiences. This could also support development of identity connotations associated with different firm types, which may appeal to workers with different cultural and social backgrounds. However, when workers are biased similarly, the presence of multiple equilibria is crucial to make it possible for workers to be able to coordinate on applying to different firms. In which case, long-term experiences have to be less prominent as compared to expected payoffs.

4 Discussions and Conclusion

In this paper, I explore the role of experiences on workers' adaptive learning behaviour in application choices. By integrating experience-based learning into search, I am able to track the evolution of firms' wages and workers' choice probabilities over time, which provides some insights on labour market dynamics. It also answers to the question of equilibrium selection and whether workers learn to apply to jobs more efficiently, defined as higher likelihood of one-to-one matching. Furthermore, this could also help to explain the apparent puzzling phenomenon of workers' lack of switching in applying to higher wage jobs even when they are able to transit (Archer (2016)), and provide an argument for sorting at application stage that is less related to skills (Barbulescu and Bidwell (2013)).

4.1 Equilibrium Selection

The presence of equilibrium multiplicity in search problem proposed, such as in this paper, often call into question which equilibrium would the workers coordinate on, and if in the long run, they would learn to apply efficiently to different firms.

For the first market structure, where workers do not observe wages and simply learn to apply based on past feedback. In fixed wage environment with multiple equilibria, workers would converge to coordinate on applying to different firms in the long run. Such learning mechanism emphasizes heavily on initial propensities and initial experiences, which can have prominent impact in determining the equilibrium chosen. As a result, if workers start off with bias towards different firms, these create a natural "lock-in" effect and is illuminate of the eventual equilibrium outcome. Similarly, positive reinforcements in the initial application rounds could propel workers into different directions and they would be "stuck" applying to the firm they are more familiar with. On the other hand, in an environment where wages are dynamically changing, wages are pushed down to 0 in the long run, and a continuum of equilibria emerge. Both workers and firms stop learning in the absence of rewards, thus the eventual equilibrium outcome simply relies on the learning path. The learning process could stop at a point where workers' have more randomized choice probabilities, and thus one-to-one matching may not be achieved.

In the second market structure, where workers observe the wages before making an application decision, their choice probabilities would depend on wages and beliefs about their opponent's choices. Since firms are setting wages based on their inference about workers' beliefs given historical realized actions, this could also constitute as a dynamic wage environment. Given multiple equilibria and that workers are highly sensitive to expected payoffs, workers could converge to the asymmetric equilibria, where they apply with high probability to different firms. For moderate wage sensitivity, however, there could be convergence to symmetric equilibria. The presence of long-term experiences as an exogenous bias could affect equilibrium multiplicity, depending on the trade-off between the magnitude of expected payoff difference and relative bias, the equilibrium selection, as a result, would be affected depending on the extent of relative bias. However, in presence of strong and diverse experience bias, there could be a unique equilibrium where workers apply more to different firms.

In general, convergence of workers' strategies to coordinating on applying to different firms can happen with experience-based learning without intervention. However, it may not be clear which asymmetric equilibria will be selected. This can be path dependent.

4.2 Mismatch Problems

These learning models portrays an evolutionary narrative which implies prolonged periods of mismatch before reaching equilibrium. In the first market structure with fixed wages, even

though the workers eventually learn to coordinate in the long run, there could be many periods of mismatch. Particularly, if workers have similar experiences, this process will be longer. Policies could cater to reducing the impact of past experiences by inducing forgetfulness. For example, filling in past working experiences and reusing them for all future applications may induce less emphasis on past experiences. As for dynamic wage environment, workers could be playing different games depending on the evolution of payoffs. It could be harder to mitigate mismatch if workers play games with a single NE, where they apply to the same firm, for too long. Policies could similarly induce forgetfulness or focus on maintaining a wage environment with pure NEs of workers applying more to different firms.

For the second market structure, given bounded beliefs about opponent behaviour, even when there is convergence to the asymmetric equilibria, there will be positive probability of overcrowding at a single firm, thus leading to mismatch problems. On top of this, when workers are less sensitive to expected payoffs, it is more likely for them to converge to applying symmetrically in the long run, adding to the likelihood of mismatch. If they are more sensitive, they will experience some periods of mismatch before converging to applying with higher probability to different firms. The extent of mismatch could be lower in this case. Adding another layer, when there exist exogenous experience bias, strong experience bias towards the same firm would induce prolonged periods of mismatch, even in the equilibrium. For weaker experiences relative to expected payoffs, in presence of multiple equilibria, could lead to similar results as that without bias. From the policy standpoint, there could be measures to reduce bias from long-term experiences, as well as policies to induce higher wage sensitivity.

4.3 Extensions

One area of extension is to incorporate learning with imprecision to portray a more realistic job search scenario. While the learning rules in the paper already contain some aspects of constrained decision-making, such as the absence of information if one does not explore options enough under reinforcement learning mechanism; and in the case of learning with best response, the inability to observe opponent's strategy and has to rely on belief updating based on past realized actions. The prospect of diving deeper into imprecision in learning the wage environment as well as noisy recalling could spark more research in cognitively-founded learning dynamics. For instance, the sequential learning environment provides basis for studying inattention. The choice of logit model could encompass noise from both inattention (Matějka and McKay (2015), Mattsson and Weibull (2002)) and experience updating, postulating a novel form of learning rule. Furthermore, even though I considered decay in memories, it could be more realistic to account for noisy memories (Da Silveira et al. (2020), Aridor et al. (2024)), which could lead to more nuanced learning mechanism, providing explanations for potential fluctuations in matching.

Last but not least, this paper provides a baseline for incorporating adaptive learning behaviour into labour market dynamics. It facilitates understanding of pairwise learning and potential equilibrium outcomes under different learning rules. An obvious extension is to scale up the model to include more firms and workers, turning it to a $I \times J$ game, and analyzing such multi-agent interactions may encompass richer results.

References

- Albarracin, D. and Wyer Jr, R. S. (2000). The cognitive impact of past behavior: influences on beliefs, attitudes, and future behavioral decisions. *Journal of personality and social psychology*, 79(1):5.
- Archer, S. . T. G. (31 May 2016). Why we don’t change jobs enough – and why we should. <https://www.theguardian.com/careers/2016/may/31/change-jobs-risks-dead-end-take-career-happiness>. Accessed: 15 Jun 2024.
- Aridor, G., da Silveira, R. A., and Woodford, M. (2024). Information-constrained coordination of economic behavior. Technical report, National Bureau of Economic Research.
- Barbulescu, R. and Bidwell, M. (2013). Do women choose different jobs from men? mechanisms of application segregation in the market for managerial workers. *Organization Science*, 24(3):737–756.
- Beggs, A. W. (2005). On the convergence of reinforcement learning. *Journal of economic theory*, 122(1):1–36.
- Briñol, P. and Petty, R. E. (2022). Self-validation theory: An integrative framework for understanding when thoughts become consequential. *Psychological Review*, 129(2):340.
- Burdett, K. and Mortensen, D. T. (1998). Wage differentials, employer size, and unemployment. *International economic review*, pages 257–273.
- Camerer, C. and Hua Ho, T. (1999). Experience-weighted attraction learning in normal form games. *Econometrica*, 67(4):827–874.
- Da Silveira, R. A., Sung, Y., and Woodford, M. (2020). Optimally imprecise memory and biased forecasts. Technical report, National Bureau of Economic Research.
- Duffy, J. and Hopkins, E. (2005). Learning, information, and sorting in market entry games: theory and evidence. *Games and Economic behavior*, 51(1):31–62.
- Eeckhout, J. (2018). Sorting in the labor market. *Annual Review of Economics*, 10(1):1–29.
- Erev, I., Ert, E., Roth, A. E., Haruvy, E., Herzog, S. M., Hau, R., Hertwig, R., Stewart, T., West, R., and Lebiere, C. (2010). A choice prediction competition: Choices from experience and from description. *Journal of Behavioral Decision Making*, 23(1):15–47.
- Erev, I. and Roth, A. E. (1998). Predicting how people play games: Reinforcement learning in experimental games with unique, mixed strategy equilibria. *American economic review*, pages 848–881.
- Fouarge, D., Kriechel, B., and Dohmen, T. (2014). Occupational sorting of school graduates: The role of economic preferences. *Journal of Economic Behavior & Organization*, 106:335–351.
- Fudenberg, D. and Levine, D. K. (1998). *The theory of learning in games*, volume 2. MIT press.
- Galenianos, M. and Kircher, P. (2009). Directed search with multiple job applications. *Journal of economic theory*, 144(2):445–471.
- Hopkins, E. (1999). A note on best response dynamics. *Games and Economic Behavior*, 29(1-2):138–150.

- Hopkins, E. (2002). Two competing models of how people learn in games. *Econometrica*, 70(6):2141–2166.
- Hopkins, E. (2007). Adaptive learning models of consumer behavior. *Journal of economic behavior & organization*, 64(3-4):348–368.
- Hopkins, E. and Posch, M. (2005). Attainability of boundary points under reinforcement learning. *Games and Economic Behavior*, 53(1):110–125.
- Kanfer, R. and Bufton, G. M. (2018). Job loss and job search: A social-cognitive and self-regulation perspective. *The Oxford handbook of job loss and job search*, pages 143–158.
- Langenhove, L. v. and Harré, R. (1994). Cultural stereotypes and positioning theory. *Journal for the Theory of Social Behaviour*, 24(4):359–372.
- Lieder, F., Griffiths, T. L., M. Huys, Q. J., and Goodman, N. D. (2018). The anchoring bias reflects rational use of cognitive resources. *Psychonomic bulletin & review*, 25:322–349.
- Lu, S. (2024). Attention as a scarce resource: Costly job search with inattentive workers. *Working Paper (Under Revision)*.
- Matějka, F. and McKay, A. (2015). Rational inattention to discrete choices: A new foundation for the multinomial logit model. *American Economic Review*, 105(1):272–298.
- Mattsson, L.-G. and Weibull, J. W. (2002). Probabilistic choice and procedurally bounded rationality. *Games and Economic Behavior*, 41(1):61–78.
- McCall, J. J. (1970). Economics of information and job search. *The Quarterly Journal of Economics*, 84(1):113–126.
- McKelvey, R. D. and Palfrey, T. R. (1995). Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1):6–38.
- Moen, E. R. (1997). Competitive search equilibrium. *Journal of political Economy*, 105(2):385–411.
- Nowé, A., Vrancx, P., and De Hauwere, Y.-M. (2012). Game theory and multi-agent reinforcement learning. *Reinforcement Learning: State-of-the-Art*, pages 441–470.
- Paas, F. and Ayres, P. (2014). Cognitive load theory: A broader view on the role of memory in learning and education. *Educational Psychology Review*, 26:191–195.
- Roth, A. E. and Erev, I. (1995). Learning in extensive-form games: Experimental data and simple dynamic models in the intermediate term. *Games and economic behavior*, 8(1):164–212.
- Terjesen, S., Vinnicombe, S., and Freeman, C. (2007). Attracting generation y graduates: Organisational attributes, likelihood to apply and sex differences. *Career development international*, 12(6):504–522.
- Vafa, K., Palikot, E., Du, T., Kanodia, A., Athey, S., and Blei, D. M. (2022). Career: A foundation model for labor sequence data. *arXiv preprint arXiv:2202.08370*.
- Van Huyck, J. B., Battalio, R. C., and Rankin, F. W. (1997). On the origin of convention: Evidence from coordination games. *The Economic Journal*, 107(442):576–596.
- van Strien, S. (2022). Dynamics of learning and iterated games lecture notes, math60007/70007/97069. Accessed: May 20, 2025.

- Wanberg, C. R., Ali, A. A., and Csillag, B. (2020). Job seeking: The process and experience of looking for a job. *Annual Review of Organizational Psychology and Organizational Behavior*, 7(1):315–337.
- Wright, R., Kircher, P., Julien, B., and Guerrieri, V. (2021). Directed search and competitive search equilibrium: A guided tour. *Journal of Economic Literature*, 59(1):90–148.
- Wu, L. (2020). *Partially Directed Search in the Labor Market*. PhD thesis, The University of Chicago.

A Proofs

A.1 Proposition 4

To show delayed adaptation in equilibrium switching, I show the following exemplary learning mechanism for 3 periods:

Example 3. *Period 1: Initial wages are randomly picked. Suppose workers start off from G1 (Figure 8), where $w_{10} > 2w_{20}$ and chose $(F1, F1)$:*

$$\text{Workers' propensities: } q_{11}^i = q_{10}^i + \frac{w_{10}}{2}, \quad q_{11}^{-i} = q_{10}^{-i} + \frac{w_{10}}{2}$$

$$\text{Firms' propensities: } \theta_{(w_{11})1}^j = \theta_{(w_{10})0}^j + (z_1 - w_{10}), \quad \theta_{(w_{21})1}^{-j} = \theta_{(w_{20})0}^{-j}$$

Workers' choice of firm 1 is reinforced; Firm 1's choice of w_{10} is reinforced, and firm 2 continue to experiment wages randomly, assuming uniform probability over its action space.

Period 2: Suppose workers continue to choose $(F1, F1)$ and firm 1 picked a new wage that is lower than the previous period, $w_{11} < w_{10}$, but the relationship $w_{11} > 2w_{21}$ still hold, then:

$$q_{12}^i = q_{11}^i + \frac{w_{11}}{2}, \quad q_{12}^{-i} = q_{11}^{-i} + \frac{w_{11}}{2}$$

$$\theta_{(w_{11})2}^j = \theta_{(w_{11})1}^j + (z_1 - w_{11}), \quad \theta_{(w_{21})2}^{-j} = \theta_{(w_{21})1}^{-j}$$

Workers' choice of firm 1 is again reinforced; Firm 1's choice of w_{11} is reinforced, and the strength of reinforcement is higher than that of a wage value equals to w_{10} . This logic implies there will be higher chance of picking lower wage values as more iterations occur. Since firm 2 was not selected in round 2, it continues to experiment within its action space randomly in the next period.

Period 3: Suppose firm 1 picked an even lower wage than the previous period, $w_{12} < w_{11}$, and the wage condition becomes $2w_{12} > w_{22} > \frac{w_{12}}{2}$. There is a switch in the game played. Given previous reinforcement, there is higher probability of selecting F1, if $(F1, F1)$ is chosen again:

$$q_{13}^i = q_{12}^i + \frac{w_{12}}{2}, \quad q_{13}^{-i} = q_{12}^{-i} + \frac{w_{12}}{2}$$

$$\theta_{(w_{12})3}^j = \theta_{(w_{12})2}^j + (z_1 - w_{12}), \quad \theta_{(w_{22})3}^{-j} = \theta_{(w_{22})2}^{-j}$$

Workers' propensities to firm 1 are positively reinforced, but with even less strength than before; firm 1's choice of lower wage is reinforced, while firm 2 continue to sample the action space randomly.

This shows that as w_{1t} is driven downwards, $w_{1t} > 2w_{2t}$ could break down, and $2w_{1t} > w_{2t} > \frac{w_{1t}}{2}$ may arise, leading to a game change from G1 to G2 (Figure 9). Choosing F1 will lead to lower reinforcement as compared to choosing F2, propensities could thus be updated, and slowly, there will be convergence towards $(F1, F2)$ and $(F2, F1)$. Nonetheless, it is possible to have multiple switching (i.e. from G1 to G2 to G3, etc.) depending on the changes in wage conditions.

Return to Section 2.2.

A.2 Proposition 5

When wage regime change from G1 to G2 at $t = \tilde{t}$, workers' start experiencing decay for propensities accumulated in G1:

$$\text{Worker } i: q_{j\tilde{t}}^i = (1 - \eta)q_{j(\tilde{t}-1)}^i, \text{ Worker } -i: q_{j\tilde{t}}^{-i} = (1 - \eta)q_{j(\tilde{t}-1)}^{-i} \quad (110)$$

And future propensities will be accumulated based on $G2$ payoffs:

$$\text{Worker } i: q_{j(\tilde{t}+1)}^i = (1 - \eta)q_{j(\tilde{t})}^i + \pi_{\tilde{t}}^i(a_{\tilde{t}}^i, a_{\tilde{t}}^{-i}, a_{\tilde{t}}^j, a_{\tilde{t}}^{-j}) \quad (111)$$

$$\text{Worker } -i: q_{j(\tilde{t}+1)}^{-i} = (1 - \eta)q_{j(\tilde{t})}^{-i} + \pi_{\tilde{t}}^{-i}(a_{\tilde{t}}^i, a_{\tilde{t}}^{-i}, a_{\tilde{t}}^j, a_{\tilde{t}}^{-j}) \quad (112)$$

Suppose payoffs in $G2$ are constant for worker i , $\pi_{\tilde{t}}^{i*}$, steady state propensities would be:

$$q_j^{i*} = \frac{\pi_j^i}{\eta} \quad (113)$$

The propensities at time $t > \tilde{t}$, where at \tilde{t} , one start new accumulation of propensity according to $G2$:

$$q_{jt}^i = (1 - \eta)^{t-\tilde{t}}q_{j\tilde{t}}^i + \sum_{k=0}^{t-\tilde{t}-1} (1 - \eta)^k \pi_j^i \quad (114)$$

$$q_{jt}^i = (1 - \eta)^{t-\tilde{t}}q_{j\tilde{t}}^i + \frac{1 - (1 - \eta)^{t-\tilde{t}}}{\eta} \pi_j^i \quad (115)$$

$$q_{jt}^i = (1 - \eta)^{t-\tilde{t}}(q_{j\tilde{t}}^i - q_j^{i*}) + q_j^{i*} \quad (116)$$

For $t > \tilde{t}$, computing propensities starting from 0, set $q_{j\tilde{t}}^i = 0$:

$$q_{jt}^i = (1 - \eta)^{t-\tilde{t}}(-q_j^{i*}) + q_j^{i*} \quad (117)$$

The time taken for propensity after \tilde{t} to exceed the ones before can be found by:

$$(1 - \eta)^{t-\tilde{t}}(-q_j^{i*}) + q_j^{i*} \geq q_{j\tilde{t}}^i(1 - \eta)^{t-\tilde{t}} \quad (118)$$

where $q_{j\tilde{t}}^i = (1 - \eta)^{\tilde{t}}q_{j0}^i + \sum_{k=0}^{\tilde{t}-1} (1 - \eta)^{\tilde{t}-k-1} \pi_k^i$.

I am interested in $t - \tilde{t}$:

$$t - \tilde{t} \leq \frac{\ln(\frac{q_j^{i*}}{q_{j\tilde{t}}^i + q_j^{i*}})}{\ln(1 - \eta)} \quad (119)$$

For small η , $\eta \in (0, 1)$, $\ln(1 - \eta) \approx -\eta$,

$$t - \tilde{t} \propto \frac{1}{\eta} \quad (120)$$

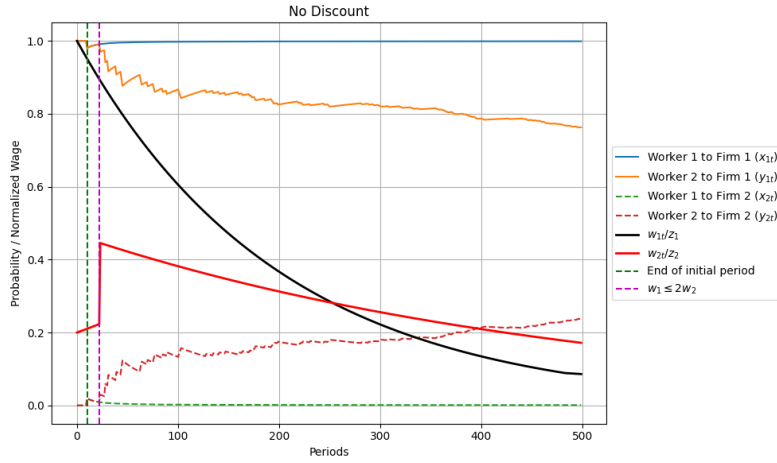
But since $G2$ payoffs are likely to be non-constant as payoffs evolve through reinforcement learning on the firms' side, the dependency of time lag on payoff changes can be more generically written as dependency on η and $f(G2 \text{ payoff dynamics})$, which captures the learning process. It could lead to longer or shorter transition process.

Return to Section 2.2.1.

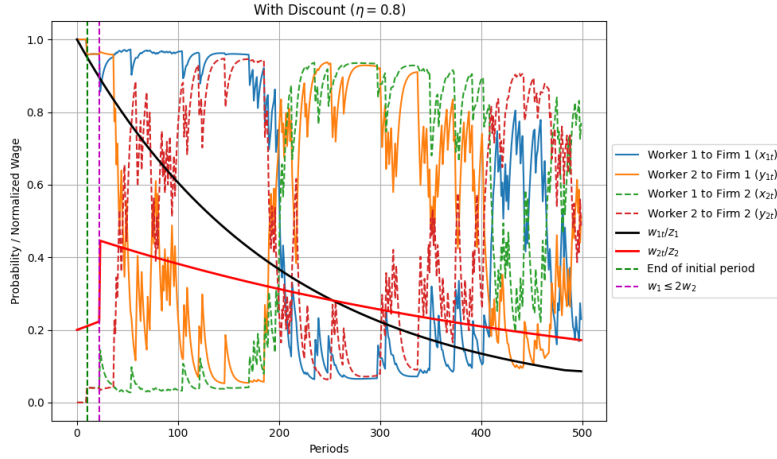
B Simulations

B.1 Example for Partial Recall of Experiences in 2-sided RL

Figure 20a and 20b show simulation of potential learning trajectory for workers with perfect and partial recall when wages are fixed to be some exogenous values that decline at a constant rate. Suppose I start with $G1$ (i.e. $w_1 \geq 2w_2$), and workers are fixed to be choosing firm 1 for 10 periods, during which, w_1 is arbitrarily adjusted downwards and w_2 upwards. Along the trajectory of exogeneous wage changes, wage condition could vary, thus shifting the game from $G1$ to $G2$ (i.e. $2w_1 > w_2 > \frac{w_1}{2}$), leading to a switch in equilibrium to learn from $(F1, F1)$ to the set $(F1, F2)$ and $(F2, F1)$ (the switching point is indicated by the pink dotted line). In this example, further shift from $G2$ to $G3$ is possible as w_1 and w_2 continues to decrease at the current rate, but not shown in the figures.



(a) Perfect Recall



(b) Partial Recall

For $z_1 = z_2 = 10$, $w_{10} = 10$, $w_{20} = 2$, all initial propensities are fixed at 1 ($x_{10} = y_{10} = 1$, $\theta_{a_0^j}^j = \theta_{a_0^j}^{-j} = 1$):

1. Assume $w_1 \geq 2w_2$ for initial 10 periods, workers are programmed to choose $(F1, F1)$, and w_1 adjust downwards by 0.5% and w_2 increases by 0.5% in each period. Over time, wage condition could reverse, and the instance where $w_1 \leq 2w_2$ is marked.
2. After 10 periods, workers are no longer fixed to choose firm 1. They can freely choose based on updated propensities and choice probabilities.
3. Once workers choose to apply to different firms, both w_1 and w_2 are programmed to decrease steadily by 0.5% per period.

Figure 20: Possible Learning Trajectory for Workers

When workers remember past events perfectly, Figure 20a shows that one of the worker would adjust his/her strategy and redirecting search towards a different firm from the other worker over time. In the long run, if the wage condition is sustained, it is expected that workers will apply with higher probability to different firms. However, the process of learning to play new equilibria can be long. There is a combined effect from learning the new set of NEs in G2, and also to overcome the initial learning experiences which direct workers to different set of NE in G1.

In Figure 20b, when there is a discount factor on past payoffs, the application strategies change more rapidly. There is a faster switch to applying more to different firms after the equilibrium switching point. But since workers are less locked-in by past experiences and put greater weight on recent payoffs, they are also more responsive to short-term positive feedback, resulting in greater fluctuations in choice probabilities. For such learning trajectory, where workers initially overcrowd at firm 1, the partial recall set-up could be beneficial in stimulating switching of job application and inducing coordination among workers at a faster pace. Workers are quicker to forget the initial experiences, thus are more adaptive to new market conditions. However, the pitfall of this is the volatility in application strategies. This could lead to a higher likelihood of mismatch, as workers' strategies are more influenced by recent events and exhibit greater stochasticity. As a result, their behaviour does not converge to applying more to different firms, even over an extended time period.

Return to Section 2.2.1.

B.2 Different Algorithms for Sequential Learning

In this section, I experimented with other algorithms, which mainly differ on the firms' side, to explore if different optimization method may affect resulting outcomes.

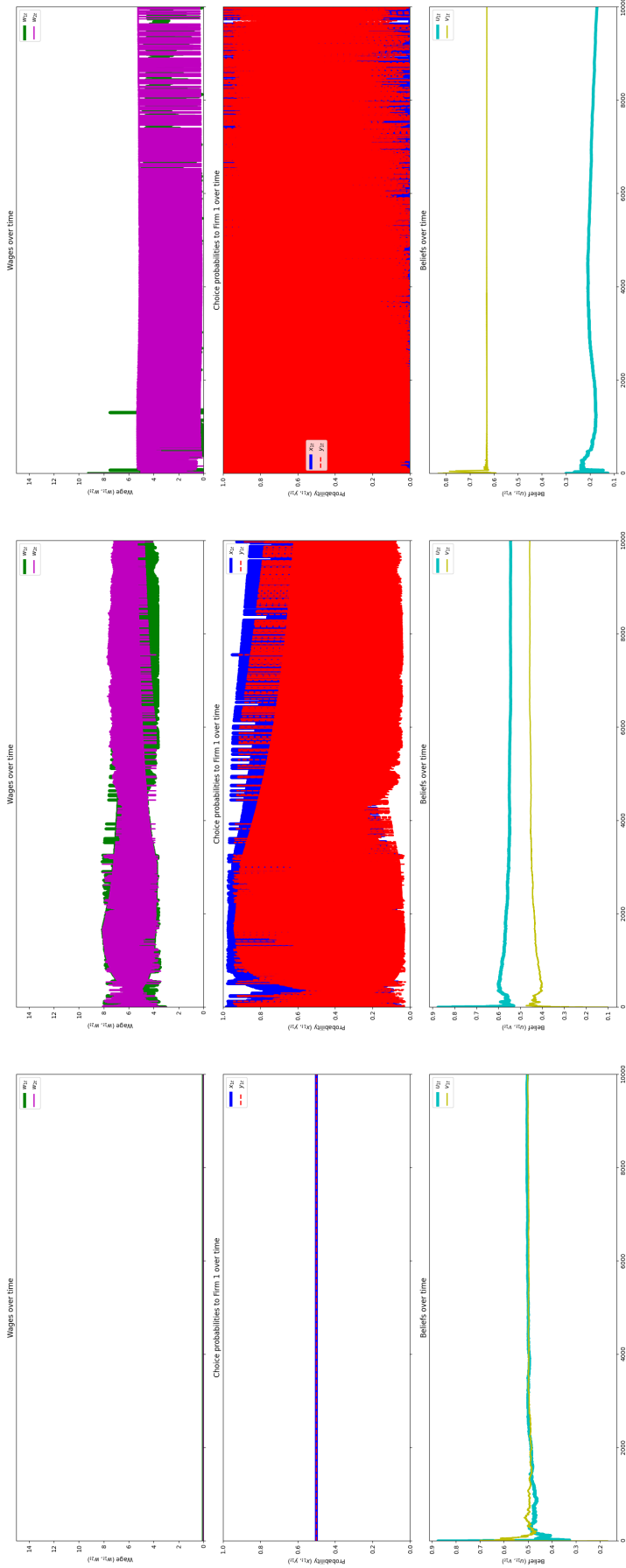
Algorithm 2 Firms' Wage-setting using Local Optimization Method 1

```

1: Initialize for  $t = 0$ , set  $\alpha_1^i$ ,  $\alpha_1^{-i}$ ,  $u_0^i$ ,  $v_0^i$ ,  $x_0$ ,  $y_0$ ; Compute  $w_{10}$ ,  $w_{20}$ .
2: for one session do
3:   Loop the following
4:   for 100000 time periods do
5:     Loop for each time period
6:     for all firms do
7:       Loop for 500 iterations
8:       for one iteration do
9:         Guess a set of wages,  $(w_{1t}^{\text{Guess}}, w_{2t}^{\text{Guess}})$ .
10:        Compute workers' reaction based on equations (41) and (42) to obtain
             $(x_{1t}^{\text{Potential}}, y_{1t}^{\text{Potential}})$ .
11:        Based on workers' potential application rates, compute wages using
            equations (54) and (55) to obtain  $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$ .
12:        Compare  $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$  with  $(w_{1t}^{\text{Guess}}, w_{2t}^{\text{Guess}})$ .
13:        If wage guess differs from the wage that would be set optimally given
            workers' response, the guess is incorrect, reiterate the process and start with
            a different guess.
14:        Make a smaller adjustment from previous guess if it was close, otherwise,
            make a bigger adjustment.
15:      end for one iteration
16:      Set wages  $(w_{1t}, w_{2t})$  to be when the difference between the guessed wages and
            potential wages is negligible, such that these are the optimal wages in each
            period given workers' reaction.
17:    end for firms
18:    for all workers do
19:      Observe  $(w_{1t}, w_{2t})$ , compute  $x_t$  and  $y_t$  given  $u_t$  and  $v_t$  using equations (41) and
            (42).
20:      Generate a choice of action from  $(F1, F2)$  for each worker based on  $x_t$  and  $y_t$ ,
            which will be the realized workers' choice in period  $t$ .
21:    end for workers
22:    for reward generation and updating do
23:      Given realized workers' actions  $(a_t^i, a_t^{-i})$  and firms' wages  $(w_{1t}, w_{2t})$ , compute
            the rewards for all agents.
24:      Workers' beliefs about each other ( $u_{t+1}$  and  $v_{t+1}$ ) are updated based on equation
            (44) for use in the following period.
25:    end for one period
26:  end for all periods
27: end for all sessions
28: return results  $(x_{1t}, x_{2t}, y_{1t}, y_{2t}, u_{1t}, u_{2t}, v_{1t}, v_{2t}, w_{1t}, w_{2t})$  for all periods.

```

With Algorithm 2, I show in Figure 21 one simulation example, which records changes in wages over time (top panel), workers' choice probability of firm 1 (middle panel) and beliefs of other



(a) $\beta = 0.2$

(b) $\beta = 1.0$

(c) $\beta = 5.0$

Figure shows wages, choice probabilities and belief evolution for different $\beta = \{0.2, 1.0, 5.0\}$, given $z_1 = z_2 = 10$, $t = 10000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0.5$, $v_{10} = 0.5$, and no smoothing.

Figure 21: Changes in Wages and Workers' Choice Probabilities Over Time

worker's choice probability to firm 1 (bottom panel). At $\beta = 0.2$, workers' choice probabilities stays at 0.5, beliefs converge to 0.5, and wages are 0 based on equation (54) and (55). At $\beta = 1.0$, wages are positive. While beliefs about opponents are seemingly stable and converge to around 0.5, wages and workers' choice probabilities fluctuate rapidly. In the final case of $\beta = 5.0$, wages are also positive. There seems to be convergence of beliefs to opponents applying with higher probability to different firms, but there is more drastic fluctuations in firms' wage setting and workers' choice probabilities.

In Table 2, I run the simulation 20 times to obtain values for 20 sessions. I then compute the average values for the last 20% of the periods for each session and listed the first and last session. I show that wages are positive when workers are more responsive to them (i.e. higher β). It is also more likely for workers to apply with higher probability to different firms when β is higher. Even though by averaging across the 20 sessions, I found for all β , the average choice probabilities concentrated on 0.5, which implies that workers are equally likely to select either of the asymmetric strategies in presence of multiple equilibria.

First and last session	w_1	w_2	x_{1t}	y_{1t}
Low β (0.2) [First]	0.000000	0.000000	0.500000	0.500000
Low β (0.2) [Last]	0.000000	0.000000	0.500000	0.500000
Medium β (1.0) [First]	4.331617	4.332060	0.499997	0.499995
Medium β (1.0) [Last]	4.331806	4.331872	0.499994	0.500004
High β (5.0) [First]	4.599850	4.600150	0.749995	0.250005
High β (5.0) [Last]	4.599290	4.600707	0.250003	0.749996
Average across sessions				
Low β (0.2)	0.000000	0.000000	0.500000	0.500000
Medium β (1.0)	4.331838	4.331839	0.500000	0.500000
High β (5.0)	5.870716	5.169464	0.462453	0.589455

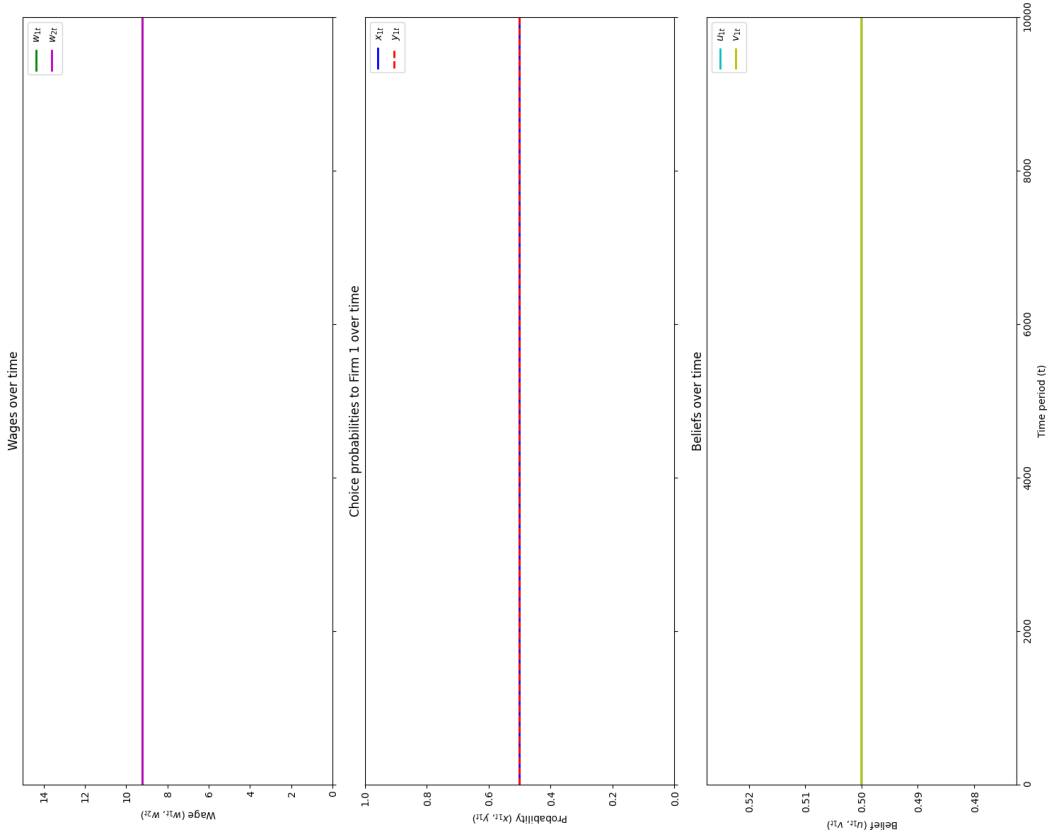
Table 2: Average Values of Last 20% of 10000 Periods for 20 Simulation Sessions

In presence of multiple equilibria when β is sufficiently large, tiny shifts in beliefs could push the system across boundaries between different basins of attractions, thereby contributing to large fluctuations in firms' and workers' behaviours. This could be the main reason for the apparent jumps across different convergence paths. However, since Algorithm 2 starts with random guesses of wages each period (Line 9), this may also contribute to large fluctuations in wages and workers' choice probabilities as values in different basin of attractions may be selected each time. This method effectively explore local optima points, where optimal wages are found near the neighbourhood of initial guesses in each period. Therefore, even with same belief inputs (u_t, v_t) , it is not guaranteed that the same convergence path is selected.

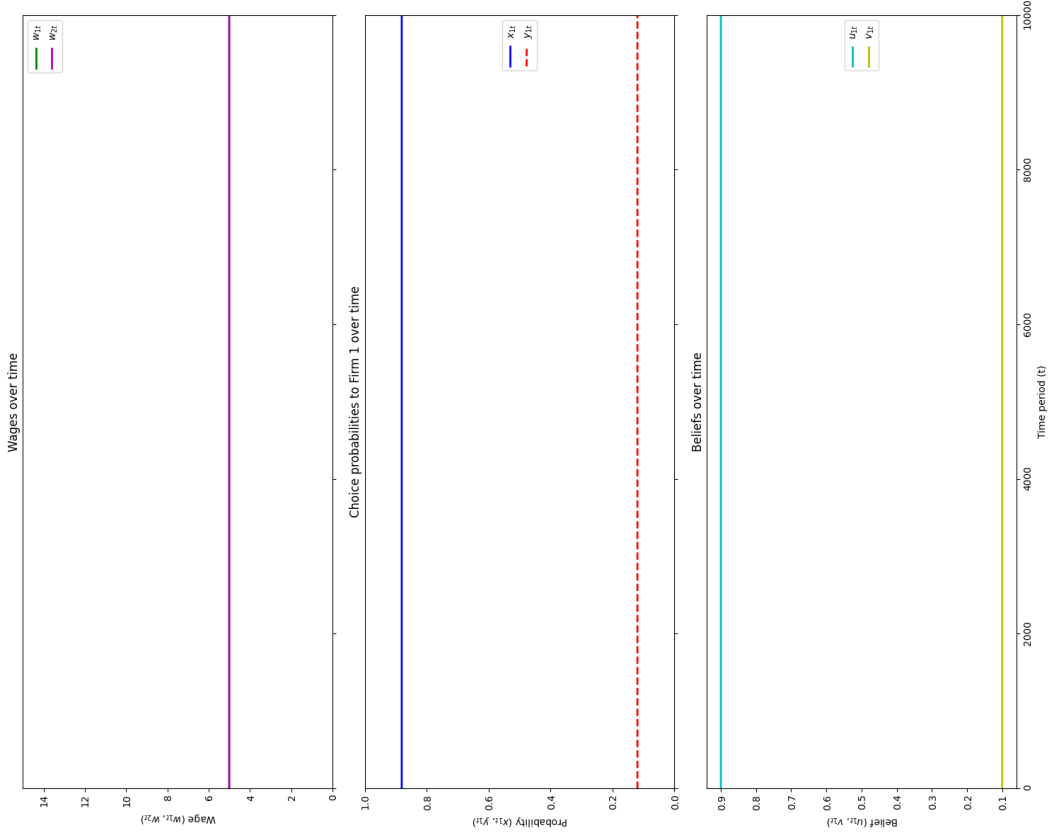
In order to boil down the reason behind the fluctuations in choice probabilities and wages, (1) I first fix workers' beliefs to investigate if workers' choice probabilities and firms' wage-setting behaviour stabilize when beliefs stabilize; and (2) I then fix both wages or only one of the wages to verify if firms' wage-setting behaviour in this algorithm is the triggering factor for the fluctuations.

Figures 22, 23 imply that if beliefs stabilize, workers and firms' behaviour would stabilize. Similarly, fixing wages would lead to convergence to a unique set of beliefs and choice probabilities. As a result, the observed fluctuations in the full system is due to small shifts in beliefs, causing fluctuations along the feedback loop:

$$u_t, v_t \rightarrow x_t, y_t \rightarrow w_{1t}, w_{2t} \rightarrow x_t, y_t \rightarrow u_t, v_t \quad (121)$$



(a) $(u_1, v_1) = (0.5, 0.5)$

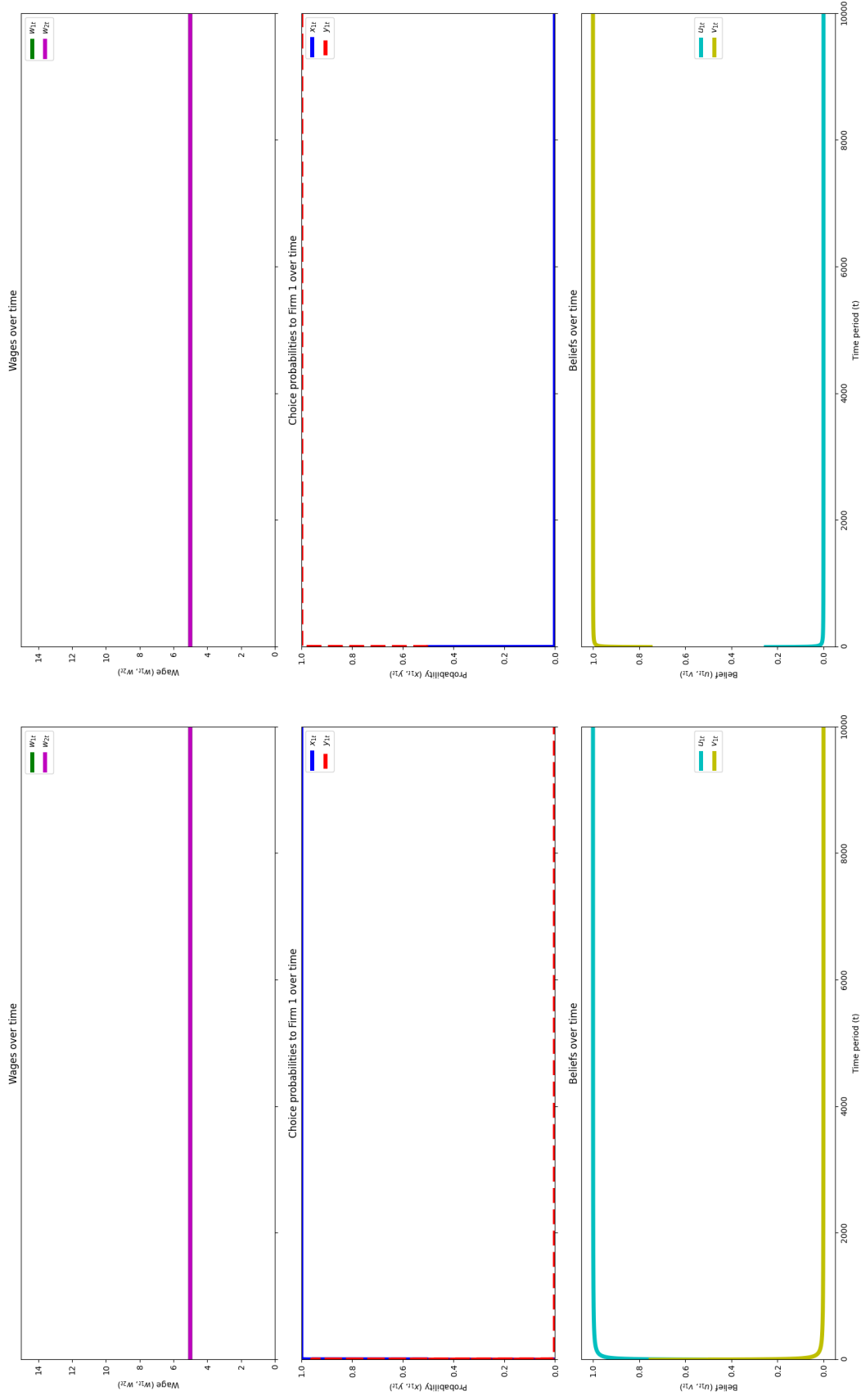


(b) $(u_1, v_1) = (0.9, 0.1)$

wages, choice probabilities, fixing beliefs at $(u_1, v_1) = (0.5, 0.5)$ and $(0.9, 0.1)$, given $\beta = 5.0$, $z_1 = z_2 = 10$, $t = 10000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0$, and no smoothing.

Figure 22: Changes in Wages, Workers' Choice Probabilities Over Time for Fixed Beliefs

Figure shows



(a) $w_1 = w_2 = 5$

(b) $w_1 = 5$

Figure shows wages, choice probabilities and belief evolution, given $\beta = 5.0$, $z_1 = z_2 = 10$, $t = 10000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0$, $u_{10} = 0.5$, $v_{10} = 0.5$, no smoothing.

Figure 23: Changes in Workers' Choice Probabilities, Beliefs Over Time for Fixed Wages

Small changes in beliefs can move workers' choice probabilities, which in turn affect wages significantly, and this further feedback into choice probabilities, leading to a magnification effect. Therefore, larger swings for choice probabilities and wages are observed in Figure 22b despite small changes in u_t and v_t .

I experiment with another algorithm (Algorithm 3), where I modify the firms' behaviours (Line 6-17) slightly to put some structure on wage guesses, such that next period wage guesses are made in the vicinity of previous period wages.

Algorithm 3 Firms' Wage-setting with Small Adjustment from Previous Period

- 1: **for** all firms **do**
 - 2: Check if previous wages are still optimal given updated workers' beliefs based on past period realized actions.
 - 3: **for** a set of $(w_{1(t-1)}, w_{2(t-1)})$ **do**
 - 4: Compute workers' $(x_{1t}^{\text{Potential}}, y_{1t}^{\text{Potential}})$ based on equations (41) and (42).
 - 5: Based on workers' potential application rates, compute corresponding wages, $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$, using equations (54) and (55).
 - 6: Compare $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$ with $(w_{1(t-1)}), (w_{2(t-1)})$.
 - 7: **end for** checking previous wages
 - 8: If previous wages is optimal, set $(w_{1t}, w_{2t}) = (w_{1(t-1)}), (w_{2(t-1)})$; adjust otherwise.
 - 9: **for** previous wages no longer optimal **do**
 - 10: Loop for 500 iterations
 - 11: **for** first iteration **do**
 - 12: Compute difference between $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$ and $(w_{1(t-1)}), (w_{2(t-1)})$.
 - 13: Make new wages guesses in small increment if wage difference is small, otherwise make a bigger adjustment:
 - 122)
$$w_{1t}^{\text{Guess}} = w_{1(t-1)} + \text{increment} * \text{wage difference} \quad (122)$$
 - 123)
$$w_{2t}^{\text{Guess}} = w_{2(t-1)} + \text{increment} * \text{wage difference} \quad (123)$$
 - 14: Given wage guesses, compute workers' reaction, $(x_{1t}^{\text{Potential}}, y_{1t}^{\text{Potential}})$, based on potential reaction, compute potential wages, $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$.
 - 15: Compare $(w_{1t}^{\text{Potential}}, w_{2t}^{\text{Potential}})$ and $(w_{1t}^{\text{Guess}}, w_{2t}^{\text{Guess}})$, if far apart, the guess is incorrect, start a new guess in subsequent iteration based on their differences.
 - 16: **end for** first iteration
 - 17: Subsequent iterations make small increments from previous guesses. The possible bound for adjustment is smaller for first 5 iterations, bigger for the next 15, and even larger afterwards. This helps to ensure convergence within 500 iterations.
 - 18: **end for** all iterations
 - 19: Set wages (w_{1t}, w_{2t}) to be when the difference between guessed wages and potential wages is negligible.
 - 20: **end for** firms
-

The simulation results are similar to Algorithm 2, again implying that small shifts in beliefs can contribute to fluctuations in choice probabilities and wages when there exist multiple equilibria.

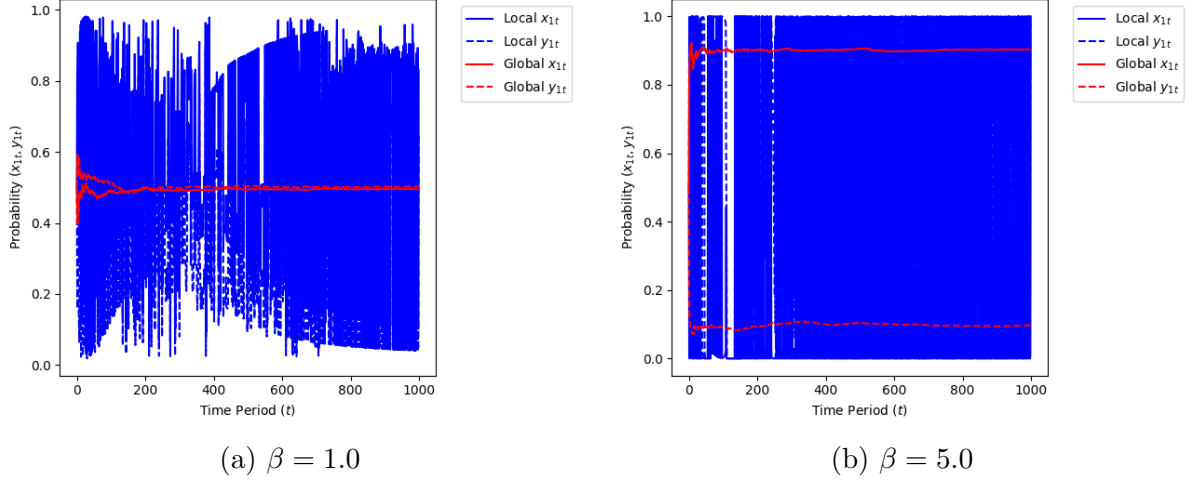


Figure shows comparison of choice probabilities evolution for the two wage-setting approach, denoted as local and global respectively. Given $\beta = \{1.0, 5.0\}$, $z_1 = z_2 = 10$, $t = 1000$, $\alpha_1^i = 0$, $\alpha_1^{-i} = 0$, $u_{10} = 0.5$, $v_{10} = 0.5$, no smoothing.

Figure 24: Comparing Workers' Strategies For the Two Wage-setting Approach

However, firms' wage-setting behaviour still plays a role. Comparing the evolution of workers' choice probabilities to firm 1 using Algorithm 1 and 3 for the cases with multiple equilibria in Figure 24, I show that there is a clearer selection of one equilibrium, and drastic fluctuations in choice probabilities and wages, respectively. The potential reason behind this is the difference between global and local optimization approach. The local optimization by Algorithm 3 can lead to circling around a single equilibrium as wages react to small adjustment in beliefs when they move from previous period values, and there can also jump across different equilibria due to large effect from the feedback loop. On the other hand, global optimization by Algorithm 1 scans the entire wage landscape for each set of beliefs, making it robust to small shifts in beliefs and avoiding such fluctuations.

Return to Section 3.1.

B.3 Workers' Side Illustration of Long-Term Experience Bias

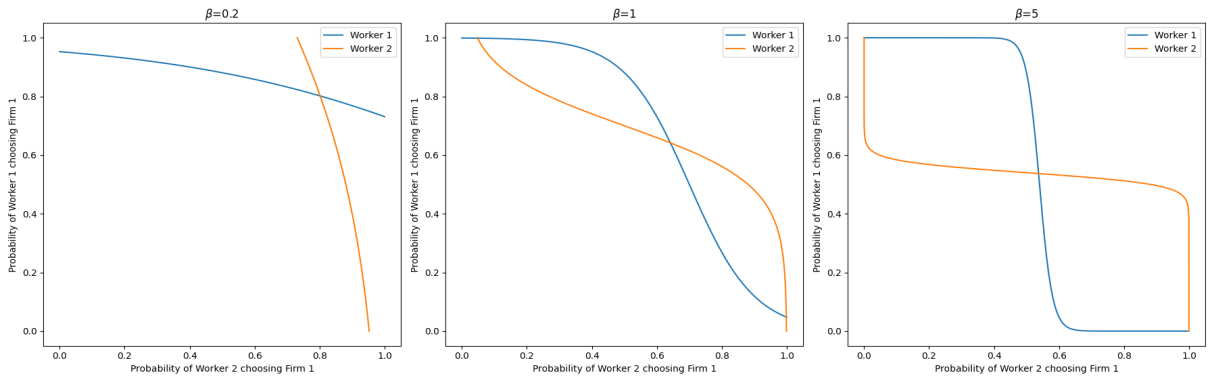


Figure show workers' choice probability for $\alpha_1^i = \alpha_1^{-i} = 2$ when fixing the wages to $w_1 = w_2 = 10$.

Figure 25: Workers' Response Functions with Bias

Return to Section 3.2