

TEMPORAL-DIFFERENCE ESTIMATION OF DYNAMIC DISCRETE CHOICE MODELS

KARUN ADUSUMILLI AND DITA ECKARDT

ABSTRACT. We propose a new algorithm to estimate the structural parameters in dynamic discrete choice models. The algorithm is based on the conditional choice probability approach, but uses the idea of Temporal-Difference learning from the Reinforcement Learning literature to estimate the different terms in the value functions. In estimating these terms with functional approximations using basis functions, our approach has the advantage of naturally allowing for continuous state spaces. Furthermore, it does not require specification of transition probabilities, and even estimation of choice probabilities can be avoided using a recursive procedure. Computationally, our algorithm only requires solving a low dimensional linear equation. We find that it is substantially faster than existing approaches when the finite dependence property does not hold, and comparable in speed to approaches that exploit this property. For the estimation of dynamic games, our procedure does not require integrating over the actions of other players, which further heightens the computational advantage. We show that our estimator is consistent, and efficient under discrete state spaces. In settings with continuous states, we propose easy to implement locally robust corrections in order to achieve parametric rates of convergence. Preliminary Monte Carlo simulations confirm the workings of our algorithm.

1. INTRODUCTION

Dynamic discrete choice (DDC) models are frequently used to describe the intertemporal choices of forward-looking individuals in a variety of contexts. In these models, agents maximize their expected future payoff through repeated choice amongst a set of discrete alternatives. Based on a revealed preference argument, structural estimation proceeds by using microdata on choices and outcomes to recover the underlying model parameters.¹ A key challenge in this literature is the complexity of estimation. Uncovering the structural parameters typically requires an explicit solution to the dynamic programming problem in addition to the optimization of an estimation criterion. In a seminal contribution, Rust (1987) develops an iterative solution algorithm, the Nested Fixed Point algorithm, that repeatedly solves the dynamic programming problem and searches for the root of the likelihood equations to update the structural parameters. To ease the computational burden associated with fully solving the dynamic optimization problem in each iteration, alternative methods have been developed. A key advance has been Hotz and Miller’s (1993) Conditional Choice Probability (CCP) algorithm which avoids the repeated solution of the intertemporal optimization problem by taking advantage of a mapping between value function differences and conditional choice probabilities. This idea has subsequently been refined by Hotz et al. (1994) who suggest a simulation-based CCP method, and Aguirregabiria and Mira (2002) who develop a more efficient recursive CCP algorithm, the nested pseudo-likelihood (NPL) algorithm. More recently, Arcidiacono and Miller (2011) exploit the property of finite dependence to speed up CCP estimation. This idea has been extended by Chernozhukov et al. (2018) to high dimensional states, also under finite dependence. Separately, Semenova (2018) also allows for high-dimensional states, but the parameters are only partially identified.

Despite these advances, the estimation of DDC models remains constrained by its computational complexity, particularly in the large class of models where finite dependence does not hold. While the CCP algorithm substantially reduces the computational burden compared to traditional methods in such settings, it becomes computationally infeasible if the number of discrete state variables is large. This problem is even more apparent when the underlying state variables are continuous and the resulting discretization gives rise to a very high-dimensional state space. An application that is particularly affected by this issue is the estimation of dynamic discrete games, where the strategic interaction of agents means that the state space increases exponentially with the number of players. Furthermore, it is uncommon for finite dependence to hold under dynamic games. Existing methods in discrete state space settings such as the pseudo-likelihood estimator proposed by Aguirregabiria and Mira (2007) or the minimum distance estimator suggested by Pesendorfer and Schmidt-Dengler (2008) become computationally difficult when the state space is large. If the states are continuous, discretization may be avoided by using forward Monte Carlo simulations (Bajari et al., 2007), but this may become very involved as the number of continuous state variables or players increases.

To overcome these limitations, we propose a new algorithm for the estimation of DDC models. Our approach is based on traditional CCP methods, but makes use of a Temporal-Difference

¹See Aguirregabiria and Mira (2010) for a detailed survey of the literature on the estimation of DDC models.

(TD) method from the Reinforcement Learning literature to provide functional approximations for the different terms in the value functions.² We start by choosing a set of basis functions in actions and observed state variables. We then project the value function operator onto the linear span of these basis functions and compute the resulting fixed point (of the projected value function operator). This fixed point is our functional approximation to the value function. Unlike most existing estimation approaches, our algorithm does not require any specification or estimation of transition probabilities. Estimating the parameters requires solving a single, low-dimensional linear equation. In the unlikely case where the dimensionality of the state space and therefore matrix makes the inversion computationally difficult, we propose an alternative stochastic gradient procedure to obtain the functional approximations for the terms in the value function. With these at hand, estimation of the structural parameters can proceed with standard methods such as maximum likelihood estimation (MLE) or minimum distance estimation.

In order to implement our functional approximation approach, an estimate for the conditional choice probabilities is required. Aguirregabiria and Mira (2002) show that, if the state variables are discrete, the error from this first-stage estimation does not have a first-order impact on the estimation of structural parameters, but this result does not carry over to a case with continuous states. Moreover, if the state variables are continuous the implicit estimation of transition probabilities in the TD algorithm has an impact on the second stage of the estimation and the pseudo maximum likelihood estimator for the structural parameters will no longer be optimal. We explain how, following Akerberg et al. (2014), Newey (1994) and Chernozhukov et al. (2018), our estimation approach for the functional approximations and the structural parameters can be easily adapted to continuous state spaces using a correction term to provide locally robust estimators. Based on Chernozhukov et al. (2018), a cross-fitting procedure is suggested to further correct for any finite sample bias resulting from the first-stage estimation. We also propose a recursive version of our algorithm, similar to the NPL algorithm by Aguirregabiria and Mira (2002), in which the conditional choice probabilities are updated as part of the estimation of the functional form approximations. Finally, we incorporate permanent unobserved heterogeneity into our methods by combining the TD estimation with an Expectation-Maximisation (EM) algorithm (Dempster et al., 1977).

We show that our estimator is consistent and converges at parametric rates. Moreover, it is computationally very cheap and therefore fast, especially in models that do not exhibit a finite dependence property. A Monte Carlo study based on the Rust (1987) bus engine replacement problem confirms the workings of our algorithm. Most importantly, our TD estimator provides a feasible estimation method when the state variables are continuous or the state space is large. This is particularly important for the estimation of dynamic discrete games. Even with discrete states, existing methods for estimation of dynamic games ((Bajari et al., 2007); Aguirregabiria and Mira (2007); Pesendorfer and Schmidt-Dengler (2008)) require integrating out the actions of the other players. With many players, or under continuous states this can get quite cumbersome. By contrast, our procedure works directly with the joint empirical distribution of the states and their sample successors. Thus the ‘integrating out’ is done implicitly

²See Sutton and Barto (2018) for details on TD learning.

within the sample expectations. In fact our estimation procedure treats single and multiple agent dynamic models in exactly the same way. The only difference is that the basis space is generally larger under dynamic games.

While most of the computational gain is achieved in models with high-dimensional state space, our approach is also as efficient as other methods in models with fewer state variables. In fact, we show that in cases where the underlying states and actions are discrete, the basis function in our functional approximations can be chosen such that our estimate is numerically identical to the one obtained from standard CCP estimators. We therefore view our method as broadening the class of DDC models that can be structurally estimated, while being as efficient as existing estimation approaches for simpler versions of the DDC problem.

In making use of a TD step in the estimation, our method relates to the literature on Reinforcement Learning. Reinforcement Learning is an area of machine learning which describes learning about how to map states into actions so as to maximize an expected payoff.³ A central component in Reinforcement Learning is the estimation of value functions. Unlike traditional dynamic programming methods, TD learning updates the current value function using sample successors. In contrast to other sample updating methods, it uses an estimate of the return instead of the actual return as target. Finally, it also employs functional approximations to approximate the value functions under continuous states. The combination of functional approximation, sample successors and estimated returns makes TD estimation extremely fast. For this reason TD algorithms are the standard method of choice for approximating value functions in Reinforcement Learning. The idea of TD learning has a long history, but the formulation in its current form is due to Sutton (1988). Tsitsiklis and van Roy (1997) studied the theoretical properties of the algorithm under functional approximation. However these were derived in the setup of online learning, whereas we intend to use our TD algorithm on a given set of observational data, i.e in an offline manner. Consequently we develop the statistical properties of TD estimation using offline data. We find that TD learning behaves very similarly to the usual series approximation in terms of convergence rates. Indeed, due to the similarity in statistical and computational properties, we like to think of TD estimation as the counterpart of least-squares regression, but for approximating value functions.

Our methods also contribute to the literature on approximating value functions. A number of techniques have been proposed for this in Economics, including parametric policy iteration (Hall et al., 2000), simulation and interpolation (Keane and Wolpin 1994), and sieve value function iteration (Arcidiacono et al., 2013). The last of these comes closest in spirit to our own approach, as the authors propose a non-parametric approximation of the value function. The difference, however, is that Arcidiacono et al. (2013) propose minimizing the TD error in the sup norm, while we minimize the projected TD error in expectation. The latter is much easier to compute and we are also able to provide strong statistical guarantees when the choice and transition probabilities are unknown, with rates of approximation that mirror standard series estimation. We also refer to Section 11.4 of Sutton and Barto (2018) for a useful discussion on the differences between minimizing the TD error and the projected TD error.

³See Sutton and Barto (2018) for a detailed treatment of Reinforcement Learning.

The remainder of this paper is organized as follows. Section 2 outlines the setup of the DDC model and fixes notation. Section 3 describes our TD estimation method for the functional approximations of the value functions, proves its theoretical properties and describes the second-step estimation of the structural parameters under discrete and continuous state variables. Section 4 describes various extensions including a recursive version of our algorithm which avoids the initial estimation of conditional choice probabilities. Section 5 incorporates permanent unobserved heterogeneity into our algorithm. Section 6 discusses the estimation of dynamic discrete games. Section 7 provides preliminary Monte Carlo simulations for our algorithm using a version of the Rust (1987) bus engine replacement problem. Section 8 concludes.

2. SETUP

We start with a single agent DDC model. Our treatment of this uses the same notation as Aguirregabiria and Mira (2010).

We consider a model in discrete time with $t = 1, \dots, T$; $T \leq \infty$ periods and $i = 1, \dots, n$ agents. We assume that the individuals are homogeneous, relegating extensions for unobserved heterogeneity to Section 5. In each period, an agent chooses among A mutually exclusive actions, each of which is denoted by a . The payoff from the action depends on the current state x . In particular, choosing action a when the state is x gives the agent an instantaneous utility of $z(a, x)^\top \theta + e$, where $z(a, x)$ is some known vector valued function of a, x and e is an idiosyncratic error term. We denote the realization of the state of an individual i at time t by x_{it} , and her corresponding action and error terms by a_{it} and e_{it} . We shall assume that e_{it} is an iid draw from some known distribution $g_e(\cdot)$. Let (a', x') denote the one-period ahead random variables immediately following the actions and states (a, x) , where $x' \sim f_X(\cdot | a, x)$. We do not make any assumptions about f_X . The utility from future periods is discounted by β .

Agent i chooses chooses actions $\mathbf{a}_i = (a_{i1}, \dots, a_{iT})$ to sequentially maximize the discounted sum of payoffs

$$E \left[\sum_{t=1}^T \beta^t \{z(x_{it}, a_{it})^\top \theta^* + e_{it}\} \right].$$

The econometrician observes the state action pairs $(\mathbf{x}_i, \mathbf{a}_i) = \{(x_{i1}, a_{i1}), \dots, (x_{iT}, a_{iT})\}$ for all individuals, but not the idiosyncratic error terms e_{it} . Using this data, the econometrician aims to recover the structural parameters θ^* . By now, a number of different algorithms have been proposed to estimate θ^* . One such algorithm, which is very popular in the literature due to its computational simplicity, is the CCP method due to Hotz and Miller (1993). This has been subsequently refined in many ways by Hotz et al. (1994), Aguirregabiria and Mira (2002), and Arcidiacono and Miller (2011), among others.

CCP methods utilize the knowledge of the conditional choice probabilities of choosing action a given state x . We shall denote these by $P_t(a|x)$ for a given period t but shall henceforth drop the subscript t with the idea that it can be made a part of the state variable x , if needed. Denote $e(a, x)$ as expected value of the idiosyncratic error term e given that action a was chosen. Hotz and Miller (1993) show that if the distribution of e follows a Generalized Extreme Value (GEV) distribution, it is possible to express $e(a, x)$ as a function of the choice probabilities

$P(a|x)$, i.e $e(a, x) = \mathcal{G}(P(a|x))$. For concreteness we shall assume in this paper that e follows a Type I Extreme Value distribution, which is perhaps the most common choice in the literature. In this case $e(a, x) = \gamma - \ln P(a|x)$, where γ is the Euler constant. Our results extend to other GEV distributions as well, after straightforward modifications. We discuss these in the Appendix.

The standard procedure in the CCP approach is as follows: Under the given distributional assumptions, the parameters are obtained as the maximizers of the pseudo-likelihood function

$$Q(\theta) = \sum_{i=1}^n \sum_{t=1}^T \log \frac{\exp \{h(a_{it}, x_{it})^\top \theta + g(a_{it}, x_{it})\}}{\sum_a \exp \{h(a, x_{it})^\top \theta + g(a, x_{it})\}},$$

where $h(\cdot)$ and $g(\cdot)$ solve the following recursive expressions:

$$\begin{aligned} h(a, x) &= z(a, x) + \beta \sum_{x'} f_X(x'|a, x) \sum_{a'} P(a'|x) h(a', x') \\ g(a, x) &= \beta \sum_{x'} f_X(x'|a, x) \sum_{a'} P(a'|x) \{e(a', x') + g(a', x')\}. \end{aligned}$$

Note that we omit the subscripts in (a, x) to denote the random variables, as opposed to realizations (a_{it}, x_{it}) . The above assumes a discrete state space. To obtain more insight, let us convert the above equations to expectations:

$$\begin{aligned} h(a, x) &= z(a, x) + \beta \mathbb{E} [h(a', x')|a, x] \\ g(a, x) &= \beta \mathbb{E} [e(a', x') + g(a', x')|a, x], \end{aligned} \tag{2.1}$$

where $\mathbb{E}[\cdot]$ denotes the expectation over the distribution of (a', x') conditional on (a, x) . Note that \mathbb{P} is a function of the distribution F of the transition and choice probabilities given by (f_X, P) . The above formulation is also valid for continuous state spaces. Both $h(a, x)$ and $g(a, x)$ have a ‘value-function’ form, which turns out to be useful as there now exist fast algorithms for computing value functions.

Observe that $h(\cdot)$ and $g(\cdot)$ are functions of the probability distributions f_X and $P(\cdot|\cdot)$, which represent the transition and conditional choice probabilities respectively. Since these are typically unknown, one usually proceeds by first estimating these as (\hat{f}_X, \hat{P}) . Typically, \hat{f}_X is obtained by MLE based on a parametric form of $f_X(x'|a, x; \theta_f)$, while \hat{P} is estimated non-parametrically using either a blocking scheme or kernel regression. Then, given (\hat{f}_X, \hat{P}) , the values of $h(\cdot)$ and $g(\cdot)$ can be estimated by solving the recursive equation 2.1. This is done by first discretizing the state space, and then solving for $h(\cdot)$, $g(\cdot)$ in terms of $z(\cdot)$, $e(\cdot)$, using either backward induction or matrix inversion.

When the underlying state variables are really continuous, discretization effectively gives rise to a very high-dimensional state space, making estimation of $h(\cdot)$, $g(\cdot)$ computationally extremely expensive. To ameliorate this issue, Hotz et al. (1994) propose forward simulation based estimators for $h(\cdot)$ and $g(\cdot)$. Nevertheless, the computational requirements remain quite high, given that such a simulation estimate has to be carried out for every possible combination of a and x . Furthermore, the simulation errors create another source of bias in small samples. Additionally, given that all the common CCP-based methods require initial estimators of θ_f and

P , these procedures often suffer from heavy bias in small samples, as θ_f and P are estimated very imprecisely and enter non-linearly in the optimization problem for the structural parameters.

In the next section we propose an alternative algorithm for maximizing $Q(\theta)$ that allows for continuous states and does not require any knowledge about or estimation of $f_X(\cdot)$. In Section 4.3, we go further and show how we can also avoid the estimation of the choice probabilities.

Notation. We fix the following notation for the rest of the paper: Let \mathbb{P} denote the population probability distribution of (a, x, a', x') . In other words, \mathbb{P} is the relative frequency of occurrence of (a, x, a', x') in the data as $n \rightarrow \infty$. Let $\mathbb{E}[\cdot]$ denote the corresponding expectation over \mathbb{P} . We shall also define $\mathbb{E}_n[\cdot]$ as the expectation over the empirical distribution \mathbb{P}_n of (a, x, a', x') . In particular, $\mathbb{E}_n[f(a, x, a', x')] := (n(T-1))^{-1} \sum_{i=1}^n \sum_{t=1}^{T-1} f(a_{it}, x_{it}, a_{it+1}, x_{it+1})$.

Let \mathcal{F} denote the space of all square integrable functions over the domain $\mathcal{A} \times \mathcal{X}$ of (a, x) . We shall use $\mathbb{E}[\cdot]$ to define a pseudo-norm $\|\cdot\|_2$ over \mathcal{F} as $\|f\|_2 := \mathbb{E}[|f(a, x)|^2]^{1/2}$ for all $f \in \mathcal{F}$.

Finally, we use $|\cdot|$ to denote the usual Euclidean norm on a Euclidean space.

3. TEMPORAL-DIFFERENCE ESTIMATION

This section presents our TD method for estimating $h(\cdot)$ and $g(\cdot)$. Let us first start with the $h(\cdot)$ function. Our method is based on a functional approximation for $h(\cdot)$. To this end, we (approximately) parameterize this as

$$h^{(j)}(a, x) \approx \phi^{(j)}(a, x)^\top \omega^{*(j)},$$

where $\phi(a, x)$ consists of a set of basis functions over the domain (a, x) , and the superscript j represents the j th dimension of $h(\cdot)$. Here, ω^* denotes some approximation weights (more on this below). For the remainder of this paper, we shall drop the superscript j indexing the dimension of $h(\cdot)$ and proceed as if the latter, and therefore θ^* , is a scalar. However, it should be taken as implicit that all our results hold for general $h(\cdot)$, as long as each dimension is treated separately. Also, to simplify the notation, we shall denote $\phi_{it} := \phi(a_{it}, x_{it})$ and $z_{it} := z(a_{it}, x_{it})$.

For any candidate function, $f(a, x)$, for $h(a, x)$, denote the TD error by

$$\delta(a, x; f) := z(a, x) + \beta \mathbb{E}[f(a', x') | a, x] - f(a, x),$$

and the dynamic programming operator by

$$\Gamma_z[f](a, x) := z(a, x) + \beta \mathbb{E}[f(a', x') | a, x].$$

Clearly, $h(a, x)$ is the unique fixed point of $\Gamma_z[\cdot]$. However we want to approximate $h(a, x)$ with a function from the linear span, \mathcal{L}_ϕ , of $\phi(a, x)$. The difficulty with the dynamic programming operator is that in general $\Gamma_z[f] \notin \mathcal{L}_\phi$ even if $f \in \mathcal{L}_\phi$. This suggests that to find a suitable approximation for $h(a, x)$ within \mathcal{L}_ϕ , we should project the dynamic programming operator back into this space. To do so, denote by P_ϕ the projection operator into the linear span of \mathcal{L}_ϕ , i.e

$$P_\phi[f](a, x) := \phi(a, x)^\top \mathbb{E}[\phi(a, x) \phi(a, x)^\top]^{-1} \mathbb{E}[\phi(a, x) f(a, x)].$$

We then obtain our approximation $\phi(a, x)^\top \omega^*$ to $h(a, x)$ as the fixed point of the projected dynamic programming operator $P_\phi \Gamma_z[\cdot]$:

$$P_\phi \Gamma_z[\phi(a, x)^\top \omega^*] = \phi(a, x)^\top \omega^*.$$

In Lemma 1 in the Appendix, we show that this in turn is equivalent to

$$\mathbb{E} [\phi(a, x) \{z(a, x) + \beta \phi(a', x')^\top \omega^* - \phi(a, x)^\top \omega^*\}] = 0, \quad (3.1)$$

which enables us to identify ω^* as

$$\omega^* = \mathbb{E} [\phi(a, x) (\phi(a, x) - \beta \phi(a', x'))^\top]^{-1} \mathbb{E} [\phi(a, x) z(a, x)]. \quad (3.2)$$

Lemma 2 in the Appendix assures that $\mathbb{E} [\phi(a, x) (\phi(a, x) - \beta \phi(a', x'))^\top]$ is indeed non-singular as long as $\beta < 1$ and $\mathbb{E} [\phi(a, x) \phi(a, x)^\top]$ is non-singular. We will discuss the properties of $\phi(a, x)^\top \omega^*$ in the next sub-section, but let us note here that in general

$$\phi(a, x)^\top \omega^* \neq P_\phi[h(a, x)].$$

Thus $\phi(a, x)^\top \omega^*$ is not the best linear approximation of $h(a, x)$, although it comes very close, as we will see shortly.

As defined above, ω^* cannot be computed directly, since it is a function of the true expectation $\mathbb{E}[\cdot]$. We can however obtain an estimator, $\hat{\omega}$, after replacing $\mathbb{E}[\cdot]$ with $\mathbb{E}_n[\cdot]$:

$$\hat{\omega} = \mathbb{E}_n [\phi(a, x) (\phi(a, x) - \beta \phi(a', x'))^\top]^{-1} \mathbb{E}_n [\phi(a, x) z(a, x)]. \quad (3.3)$$

Using the above, we obtain an estimate of $h(\cdot)$ as $\hat{h}(a, x) = \phi(a, x)^\top \hat{\omega}$.

We now turn to the estimation of $g(\cdot)$. We approximate $g(\cdot)$ using basis functions $r(a, x)$:

$$g(a, x) \approx r(a, x)^\top \xi^*.$$

As before, denote by $f(a, x)$ any candidate function for $g(a, x)$. Note that $g(a, x)$ is the unique fixed point of the operator $\Gamma_e[\cdot]$, where

$$\Gamma_e[f](a, x) := \beta \mathbb{E}[e(a', x') + f(a', x') | a, x].$$

We then obtain our approximation $r(a, x)^\top \xi^*$ to $g(a, x)$ as the fixed point of the projected operator $P_r \Gamma_e[\cdot]$, where P_r is the projection operator into the linear span of r , i.e

$$P_r[f](a, x) := r(a, x)^\top \mathbb{E}[r(a, x) r(a, x)^\top]^{-1} \mathbb{E}[r(a, x) f(a, x)].$$

As before, we may equivalently write

$$\mathbb{E} [r(a, x) \{\beta e(a', x') + \beta r(a', x')^\top \xi^* - r(a, x)^\top \xi^*\}] = 0. \quad (3.4)$$

This allows us to identify ξ^* as

$$\xi^* = \mathbb{E} [r(a, x) (r(a, x) - \beta r(a', x'))^\top]^{-1} \mathbb{E} [\beta r(a, x) e(a', x')].$$

Assuming $e(a, x)$ is known, this suggests the following estimator for ξ^* :

$$\hat{\xi} = \mathbb{E}_n [r(a, x) (r(a, x) - \beta r(a', x'))^\top]^{-1} \mathbb{E}_n [\beta r(a, x) e(a', x')]. \quad (3.5)$$

In general the term $e(a, x) = \gamma - \ln P(a|x)$ is a function of choice probabilities, which are unknown. Thus we first need to non-parametrically estimate them. Let us denote $\eta(a, x) := P(a|x)$. Suppose that we have access to a non-parametric estimator $\hat{\eta}$ of η . This can be obtained in many ways, e.g through series or kernel regression. We can then plug in this estimate to obtain $e(a, x; \hat{\eta}) := \gamma - \ln \hat{\eta}(a, x)$. This in turn enables us to obtain $\hat{\xi}$ as

$$\hat{\xi} = \mathbb{E}_n [r(a, x) (r(a, x) - \beta r(a', x'))^\top]^{-1} \mathbb{E}_n [\beta r(a, x) e(a', x'; \hat{\eta})]. \quad (3.6)$$

Using the above, we obtain an estimate of $g(\cdot)$ as $\hat{g}(a, x) = r(a, x)^\top \hat{\xi}$.

In the discrete setting, it turns out that the estimation error from $\hat{\eta}$ is not first order relevant for the estimation of θ^* , as long as θ^* is estimated using a pseudo-MLE. This was first noted in Aguirregabiria and Mira (2002).

In fact, even with continuous states, the estimation of $\hat{\xi}$ is unaffected to a first order by the estimation of $\hat{\eta}$, even though the latter only converges to the true η at non-parametric rates. This is because an orthogonality property holds for the estimation of ξ , in that

$$\partial_\eta \mathbb{E} [\beta r(a, x) e(a', x'; \eta)] = 0, \quad (3.7)$$

where $\partial_\eta \cdot$ denotes the Fréchet derivative with respect to η . To show (3.7), let us first expand the term $\mathbb{E} [\beta r(a, x) e(a', x'; \eta)]$ as follows

$$\begin{aligned} \mathbb{E} [\beta r(a, x) e(a', x'; \eta)] &= \mathbb{E} [\beta r(a, x) \mathbb{E} [e(a', x'; \eta) | a, x, x']] \\ &= \mathbb{E} [\beta r(a, x) \mathbb{E} [e(a', x'; \eta) | x']] \\ &= \mathbb{E} [\beta r(a, x) \mathbb{E} [\gamma - \ln \eta(a', x') | x']], \end{aligned} \quad (3.8)$$

where the second equality follows from the Markov property. Now, it turns out that

$$\partial_\eta \mathbb{E} [\ln \eta(a', x') | x'] = 0.$$

Indeed, consider the expression $M(\tilde{\eta}) := \mathbb{E} [\ln \tilde{\eta}(a', x') | x']$, evaluated at different candidate values $\tilde{\eta}(\cdot, \cdot)$. When evaluated at the true conditional choice probability, i.e when $\tilde{\eta}(\cdot, \cdot) = \eta(\cdot, \cdot)$, $M(\tilde{\eta})$ becomes the conditional entropy and attains its maximum. Formally, for any candidate $\tilde{\eta}(\cdot, \cdot)$, we have

$$\mathbb{E} [\ln \tilde{\eta}(a', x') | x'] - \mathbb{E} [\ln \eta(a', x') | x'] = \mathbb{E} \left[\ln \frac{\tilde{\eta}(a', x')}{\eta(a', x')} \middle| x' \right] \leq \ln \mathbb{E} \left[\frac{\tilde{\eta}(a', x')}{\eta(a', x')} \middle| x' \right] = 0.$$

This proves $\partial_\eta \mathbb{E} [\ln \eta(a', x') | x'] = 0$. Consequently, in view of (3.8), it follows that (3.7) holds. Thus, $\hat{\xi}$ is a locally robust estimator for ξ .

Even with a locally robust estimator, the use of a non-parametric estimator may lead to substantial finite sample bias. For this reason, we advocate a cross-fitting procedure (see Chernozhukov et al., 2018). In our context, this entails the following: we randomly partition the data into two folds.⁴ We estimate $\hat{\xi}$ separately for each fold using $\hat{\eta}$ estimated from the opposite fold. The final estimate of ξ^* is the weighted average of $\hat{\xi}$ from both the folds.

⁴One could of course use any finite number of folds, though we describe the method with two folds for simplicity.

Note that computation of $\hat{\omega}$ and $\hat{\xi}$ only involves inverting a $(k \times k)$ -dimensional matrix, where k is the dimension of ϕ . This is computationally extremely cheap. Using $\hat{h}(a, x)$ and $\hat{g}(a, x)$, we can in turn estimate θ^* in many different ways. For instance, we can use the pseudo-MLE estimator

$$\hat{\theta} := \arg \max_{\theta} \hat{Q}(\theta) := \sum_{i=1}^n \sum_{t=1}^{T-1} \log \frac{\exp \left\{ \hat{h}(a_{it}, x_{it}) \theta + \hat{g}(a_{it}, x_{it}) \right\}}{\sum_a \exp \left\{ \hat{h}(a, x_{it}) \theta + \hat{g}(a, x_{it}) \right\}}. \quad (3.9)$$

It turns out the estimate from (3.9) is suboptimal under continuous states. We discuss this in greater detail in Section 3.2, where we suggest a locally robust version of (3.9).

3.1. Discrete states. Suppose that the underlying states and actions are discrete, and that our algorithm uses basis functions comprised of the set of all discrete elements of x, a . Then the resulting estimate of $h(a, x)$ obtained from our algorithm is exactly the same as that obtained from the standard CCP estimators, if both the choice and transition probabilities were estimated using cell values. To see this, we note the following: First, the standard CCP estimators (see e.g. Aguirregabiria and Mira, 2010), estimate $h(a, x)$ by solving the recursive equations

$$\check{h}(a, x) = z(a, x) + \beta \sum_{x'} \hat{f}_X(x'|a, x) \sum_{a'} \hat{P}(a'|x') \check{h}(a', x'), \quad (3.10)$$

where \hat{f}, \hat{P} are estimates of f, P obtained as cell estimates. Second, by the results of Tsitsiklis and Van Roy (1997), it can be shown that when the functional approximation saturates all the states, the TD estimate from (3.3), denoted by $\hat{h}(x, a) := \phi(a, x)^\top \hat{\omega}$ satisfies the equation

$$z(a, x) + \beta \mathbb{E}_n[\hat{h}(a', x')|a, x] = \hat{h}(a, x),$$

where $\mathbb{E}_n[\hat{h}(a', x')|a, x]$ denotes the conditional expectation of $\hat{h}(a', x')$ given a and x under the empirical distribution \mathbb{P}_n (the conditional distribution exists because of the discrete number of states). But for discrete data, $\mathbb{E}_n[\hat{h}(a', x')|a, x]$ is simply

$$\mathbb{E}_n[\hat{h}(a', x')|a, x] = \sum_{x'} \hat{f}_X(x'|a, x) \sum_{a'} \hat{P}(a'|x') \hat{h}(a', x'),$$

and the value of $\hat{h}(a, x)$ and $\check{h}(a, x)$ coincide exactly. Thus, the two algorithms give identical results (a similar property also holds for $g(a, x)$). Since our estimates $\hat{h}(a, x)$ coincide with those from the standard CCP estimators, the resulting estimate $\hat{\theta}$ is also exactly the same. As a result, the final estimates of θ from both procedures also coincide exactly.

When the states are discrete, Aguirregabiria and Mira (2002) show that the estimation of η is orthogonal to the estimation of θ^* . This holds true for our procedure as well since our estimator is numerically equivalent to the one proposed by Aguirregabiria and Mira (2002). It is important to note, however, that the estimation of the transition probabilities $f_X(x'|a, x)$ is not orthogonal to the estimation of θ^* . This is not too much of an issue with discrete states since any estimate, $\hat{f}_X(x'|a, x)$, of $f_X(x'|a, x)$ converges at parametric rates, so \sqrt{n} consistent estimation of θ is still possible. However, as we will see in Section 3.3, this creates issues once we move to continuous states.

3.2. Theoretical Properties of TD estimators under continuous states. We now characterize the formal properties of our TD fixed point estimates of $h(\cdot)$ and $g(\cdot)$. We shall only focus on the case of continuous states, since under discrete states, our procedure gives exactly the same output as previous methods.

We start by characterizing the estimation error of $h(\cdot)$. Let k_ϕ denote the dimension of ϕ . We shall take $k_\phi \rightarrow \infty$ as $n \rightarrow \infty$. We impose the following assumptions for the estimation of $h(a, x)$:

Assumption 1. (i) *The basis vector $\phi(a, x)$ is linearly independent (i.e. $\phi(a, x)^\top \omega = 0$ for all (a, x) if and only if $\omega = 0$). Additionally, the eigenvalues of $\mathbb{E}[\phi(a, x)\phi(a, x)^\top]$ are uniformly bounded away from zero for all k_ϕ .*

(ii) *The basis functions are uniformly bounded, i.e. $|\phi(a, x)|_\infty \leq M$ for some $M < \infty$.*

(iii) *There exists $C < \infty$ and $\alpha > 0$ such that $\|h(a, x) - P_\phi[h(a, x)]\|_2 \leq Ck_\phi^{-\alpha}$.*

(iv) *The domain of (a, x) is a compact set, and there exists $L < \infty$ such that $|z(a, x)|_\infty \leq L$.*

(v) *$k_\phi \rightarrow \infty$ and $k_\phi^2/n \rightarrow 0$ as $n \rightarrow \infty$.*

Assumption 1(i) rules out multi-collinearity in the basis functions. This is easily satisfied. Assumption 1(ii) ensures that the basis functions are bounded. This is again a mild requirement and is easily satisfied if either the domain of (a, x) is compact, or the basis functions are chosen appropriately (e.g a Fourier basis). Assumption 1(iii) is a standard condition on the rate of approximation of $h(a, x)$ using a basis approximation. The value of α is related to the smoothness of $h(\cdot)$. Newey (1997) shows that for splines and power series, we can set $\alpha = r/d$, where r is the number of continuous derivatives of $h(\cdot)$, and d is the dimension of (a, x) . Similar results can also be derived for other approximating functions such as Fourier series, wavelets and Bernstein polynomials. The smoothness properties of $h(a, x)$ are discussed in Appendix B, where we provide some primitive conditions on $z(a, x)$, $f_X(x'|a, x)$ that ensure existence of r continuous derivatives of $h(a, x)$. Assumption 1(iv) requires the function $z(a, x)$ to be bounded. A sufficient condition for this is that $z(a, x)$ is continuous (since its domain is bounded).

Finally, Assumption 1(v) specifies the rate at which the dimension of the basis functions are allowed to grow. The rate requirements are also mild, and are the same as those employed for standard series estimation, even though our procedure is not the same as series estimation. For the theoretical properties, the exact rate of k_ϕ is not relevant up to a first order since we propose estimators of θ^* that are locally robust to estimation of $g(\cdot)$. But the choice of k_ϕ could matter in practice. For this reason we propose selecting k_ϕ through a procedure akin to cross-validation. The value of ω is estimated using a training sample and its performance evaluated on a hold-out or test sample. However in contrast to standard cross-validation, the performance is measured in terms of the empirical MSE of the TD error $\mathbb{E}_{\text{test}}[\delta^2(a, x; \hat{g})]$ on the test dataset. The value of k_ϕ that is chosen is the one that achieves the lowest MSE on the TD error.

We then have the following theorem on the estimation of $h(a, x)$:

Theorem 1. *Under Assumptions 1(i) to 1(v), the following hold:*

(i) *Both ω^* and $\hat{\omega}$ exist, the latter with probability approaching one.*

(ii) $\|h(a, x) - \phi(a, x)^\top \omega^*\|_2 \leq (1 - \beta)^{-1} \|h(a, x) - P_\phi h(a, x)\|_2 \leq C(1 - \beta)^{-1} k_\phi^{-\alpha}$.

(iii) There exists some $C < \infty$ such that with probability approaching one,

$$|\hat{\omega} - \omega^*| \leq C \sqrt{\frac{k_\phi}{n}}.$$

(iv) The L^2 error for the difference between $h(a, x)$ and $\phi(a, x)^\top \hat{\omega}$ is bounded as

$$\|h(a, x) - \phi(a, x)^\top \hat{\omega}\|_2 = O_p \left(\frac{k_\phi}{\sqrt{n}} + k_\phi^{-\alpha} \right).$$

We prove Theorem 1 in the Appendix by adapting the results of Tsitsiklis and Van Roy (1997). The first part of Theorem 1 assures that both population and empirical TD fixed points exist. The second part of Theorem 1 implies the approximation bias from $\phi(a, x)^\top \omega^*$ is within a $(1 - \beta)^{-1}$ factor of that from $P_\phi h(a, x)$. Note that the latter is the best one could do under an L_2 norm, so the theorem assures that we are only a constant away from attaining this. The third part of Theorem 1 characterizes the rate of convergence of $\hat{\omega}$ to ω^* , and the final part of Theorem 1 characterizes the rate of estimation of $h(a, x)$ itself.

In a similar vein, we impose the following assumptions for the estimation of $g(a, x)$. Let k_r denote the dimension of $r(a, x)$.

Assumption 2. (i) The basis vector $r(a, x)$ is linearly independent, and the eigenvalues of $\mathbb{E}[r(a, x)r(a, x)^\top]$ are uniformly bounded away from zero for all k_r .

(ii) $|r(a, x)|_\infty \leq M$ for some $M < \infty$.

(iii) There exists $C < \infty$ and $\alpha > 0$ such that $\|g(a, x) - P_r[g(a, x)]\|_2 \leq C k_r^{-\alpha}$.

(iv) The domain of (a, x) is a compact set, and $|e(a, x)|_\infty \leq L < \infty$.

(v) $k_r \rightarrow \infty$ and $k_r^2/n \rightarrow 0$ as $n \rightarrow \infty$.

(vi) $\hat{\xi}$ is estimated from a cross-fitting procedure described above. The conditional choice probability function satisfies $\eta(a, x) \geq \delta > 0$, where δ is independent of a, x . Additionally, $|\eta(a, x) - \hat{\eta}(a, x)|_\infty = o_p(1)$ and $\|\eta(a, x) - \hat{\eta}(a, x)\|_2^2 = o_p(n^{-1/2})$.

Assumption 2 is a direct analogue of Assumption 1, except for the last part which provides regularity conditions when $\eta(\cdot)$ is estimated. These conditions are typical for locally robust estimates and only require the non-parametric function $\eta(a, x)$ to be estimable at faster than $n^{-1/4}$ rates. This is easily verified for most non-parametric estimation methods such as kernel or series regression. Under these assumptions, we have the following analogue of Theorem 1.

Theorem 2. Under Assumptions 2(i) to 2(v), the following hold:

(i) Both ξ^* and $\hat{\xi}$ exist, the latter with probability approaching one.

(ii) $\|g(a, x) - r(a, x)^\top \xi^*\|_2 \leq (1 - \beta)^{-1} \|g(a, x) - P_r g(a, x)\|_2 \leq C(1 - \beta)^{-1} k_r^{-\alpha}$.

(iii) There exists some $C < \infty$ such that with probability approaching one,

$$|\hat{\xi} - \xi^*| \leq C \sqrt{\frac{k_r}{n}}.$$

(iv) The L^2 error for the difference between $g(a, x)$ and $r(a, x)^\top \hat{\xi}$ is bounded as

$$\|g(a, x) - r(a, x)^\top \hat{\xi}\|_2 = O_p \left(\frac{k_r}{\sqrt{n}} + k_r^{-\alpha} \right).$$

Theorem 1 and Theorem 2 imply that we can estimate $h(a, x)$ and $g(a, x)$ at reasonably fast rates. However we still need to discuss how this relates to consistent estimation of θ^* . We do this below.

3.3. Continuous states and locally robust estimation. When the states are continuous, estimation of $h(a, x)$ and $g(a, x)$ is inherently non-parametric. Unlike the case with discrete states, the estimation error from the non-parametric functions does affect the estimation of θ^* to a first order, when using the pseudo-MLE criterion. The reason for this is that $h(a, x)$ and $g(a, x)$ are actually functions of two non-parametric terms: the choice probabilities $\eta(a, x)$, and the transition probabilities $f_X(x'|a, x)$. The TD estimator implicitly takes both into account with a series approximation. Since the estimates for $f_X(x'|a, x)$ and θ^* are not orthogonal under a pseudo-MLE, this extends to the lack of orthogonality between the estimates for $h(a, x)$, $g(a, x)$ and θ^* . Consequently, the estimator, $\hat{\theta}$, based on partial likelihood will converge at slower than parametric rates.

3.3.1. Construction of the locally robust estimator. We now describe the construction of a locally robust version of the pseudo-MLE estimator. For the present analysis, let us suppose that $h(x, a)$ and $g(x, a)$ are finite-dimensional, i.e $h(x, a) \equiv \phi(x, a)^\top \omega^*$ and $g(x, a) \equiv r(x, a)^\top \xi^*$. Denote $\mathbf{v} := (\omega, \xi)$, $\mathbf{v}^* := (\omega^*, \xi^*)$ and let

$$Q(a, x; \theta, \mathbf{v}) = \ln \pi_{\theta, \mathbf{v}}(a, x); \quad \pi_{\theta, \mathbf{v}}(a, x) := \frac{\exp \{(\phi(a, x)^\top \omega) \theta + r(a, x)^\top \xi\}}{\sum_{\check{a}} \exp \{(\phi(\check{a}, x)^\top \omega) \theta + r(\check{a}, x)^\top \xi\}}.$$

The true value θ^* is then

$$\theta^* = \arg \max_{\theta} \mathbb{E} [Q(a, x; \theta, \mathbf{v}^*)].$$

Since the criterion function is convex, we can alternatively identify θ^* using the moment function

$$\mathbb{E}[m(a, x; \theta^*, \mathbf{v}^*)] = 0; \quad m(a, x; \theta, \mathbf{v}) := \partial_{\theta} Q(a, x; \theta, \mathbf{v}). \quad (3.11)$$

The lack of orthogonality of the estimator based on (3.11) is evident by the fact $\partial_{\mathbf{v}} \mathbb{E}[m(a, x; \theta, \mathbf{v}^*)] \neq 0$. Note that ω^* and ξ^* are in turn estimated using the moment functions

$$\mathbb{E}[\varphi_h(a, x, \omega^*)] = 0, \text{ and } \mathbb{E}[\varphi_g(a, x, \xi^*)] = 0, \quad (3.12)$$

where, given (3.1) and (3.4),

$$\begin{aligned} \varphi_h(a, x, \omega) &:= \phi(a, x) z(a, x) + \phi(a, x) (\beta \phi(a', x') - \phi(a, x))^\top \omega, \text{ and} \\ \varphi_g(a, x, \xi) &:= \beta r(a, x) e(a', x'; \hat{\eta}) + r(a, x) (\beta r(a', x') - r(a, x))^\top \xi. \end{aligned}$$

We make use of (3.11) and (3.12) to construct a locally robust moment for θ^* . Following Newey (1994), Akerberg et al. (2014) and Chernozhukov et al. (2018), this is given by

$$\mathbb{E}[\zeta(a, x; \theta^*, \mathbf{v}^*)] = 0, \quad (3.13)$$

where

$$\begin{aligned} \zeta(a, x; \theta, \mathbf{v}) &:= m(a, x; \theta, \mathbf{v}) - \mathbb{E}[\partial_{\omega} m(a, x; \theta, \mathbf{v})] \mathbb{E}[\partial_{\omega} \varphi_h(a, x, \omega)]^{-1} \varphi_h(a, x, \omega) \\ &\quad - \mathbb{E}[\partial_{\xi} m(a, x; \theta, \mathbf{v})] \mathbb{E}[\partial_{\xi} \varphi_g(a, x, \xi)]^{-1} \varphi_g(a, x, \xi). \end{aligned}$$

Note that

$$\begin{aligned}\mathbb{E}[\partial_\omega \varphi_h(a, x, \omega)] &= \mathbb{E}[\phi(a, x) (\beta \phi(a', x') - \phi(a, x))^\top], \text{ and} \\ \mathbb{E}[\partial_\xi \varphi_g(a, x, \xi)] &= \mathbb{E}[r(a, x) (\beta r(a', x') - r(a, x))^\top].\end{aligned}$$

We can now construct a locally robust estimator for θ^* based on (3.13). Following Chernozhukov et al. (2018), we employ a cross-fitting procedure by randomly splitting the data into two samples \mathcal{N}_1 and \mathcal{N}_2 . We compute $\hat{\omega}$ and $\hat{\xi}$ using one of the samples, say \mathcal{N}_2 . Denote by $\mathbb{E}_n^{(1)}[\cdot]$ the empirical expectation using only the observations in the first sample. We then obtain $\hat{\theta}$ as the solution to the moment equation

$$\mathbb{E}_n^{(1)}[\zeta_n(a, x; \theta, \hat{\omega}, \hat{\xi})] = 0, \quad (3.14)$$

where

$$\begin{aligned}\zeta_n(a, x; \theta, \mathbf{v}) &:= m(a, x; \theta, \mathbf{v}) - \mathbb{E}_n^{(1)}[\partial_\omega m(a, x; \theta, \mathbf{v})] \mathbb{E}_n^{(1)}[\partial_\omega \varphi_h(a, x, \omega)]^{-1} \varphi_h(a, x, \omega) \\ &\quad - \mathbb{E}_n^{(1)}[\partial_\xi m(a, x; \theta, \mathbf{v})] \mathbb{E}_n^{(1)}[\partial_\xi \varphi_g(a, x, \xi)]^{-1} \varphi_g(a, x, \xi).\end{aligned}$$

The use of cross-fitting or sample splitting is critical. If we had used the entire sample to estimate all of θ^* , ω^* and ξ^* , we would have $\mathbb{E}_n[\varphi_g(a, x, \hat{\omega})] = 0$ for $g \in \{h, e, \eta\}$, which implies $\mathbb{E}_n[\zeta_n(a, x, \theta, \hat{\omega}, \hat{\xi})] = \mathbb{E}_n[m(a, x, \theta, \hat{\omega}, \hat{\xi})]$. As noted by Chernozhukov et al. (2018), cross-fitting gets rid of the ‘own observation bias’ that is the source of the degeneracy here.

We will refer to the solution $\hat{\theta}$ of (3.14) as the locally robust pseudo-MLE estimator of θ^* . Note that we would need three way sample splits if we employ cross-fitting procedures for both estimation of θ^* and ξ^* . But the use of cross-fitting for $\hat{\xi}$ is not as critical as that for $\hat{\theta}$, and can be avoided if necessary.

Estimation of $\hat{\theta}$ using (3.14) involves non-convex optimization. Since this could cause difficulties in practice, we recommend a two-step method for computation. We first obtain a preliminary estimate $\hat{\theta}_1$ by solving the empirical analogue of (3.11). This is a convex optimization problem, and is usually very fast. Note that $\hat{\theta}_1$ is consistent for θ , even though its not efficient. We can then use $\hat{\theta}_1$ as the starting point for a Newton-Raphson or some other gradient descent algorithm for finding the root of (3.14).

3.3.2. Non-parametric analysis. For the setup of finite-dimensional $h(a, x)$ and $g(a, x)$, it is straightforward to show that the above procedure leads to \sqrt{n} rates of estimation of θ^* (see e.g. Newey (1994)). In this paper, we are primarily interested in the case where these quantities are infinite-dimensional. Still, treating the first step as parametric leads to an estimation strategy that is also valid non-parametrically as long as we let the series terms grow to infinity. To show this, we will need to derive the exact form of the adjustment terms in the non-parametric case. To this end, we will make use of the form of the parametric adjustment terms in (3.14) to conjecture the expression for the non-parametric correction term. We shall then verify that this indeed leads to a locally robust estimator.

With the above in mind, consider the adjustment term

$$\hat{\mathcal{A}}_h := \mathbb{E}_n^{(1)}[\partial_\omega m(a, x; \theta, \mathbf{v})] \mathbb{E}_n^{(1)}[\partial_\omega \varphi_h(a, x, \omega)]^{-1} \varphi_h(a, x, \omega)$$

for $h(a, x)$. Denote

$$m(a, x; \theta, h, g) := \partial_\theta Q(a, x; \theta, h, g); \quad Q(a, x; \theta, h, g) := \ln \frac{\exp \{h(a, x)\theta + g(a, x)\}}{\sum_{\tilde{a}} \exp \{h(\tilde{a}, x)\theta + g(\tilde{a}, x)\}}.$$

Then $\hat{\mathcal{A}}_h$ can be rewritten as

$$\hat{\mathcal{A}}_h = \hat{\lambda}_h(a, x) \{z(a, x) + \beta \phi(a', x')^\top \omega - \phi(a, x)^\top \omega\}, \quad (3.15)$$

where

$$\hat{\lambda}_h(a, x) := \phi(a, x)^\top \mathbb{E}_n^{(1)} [(\beta \phi(a', x') - \phi(a, x)) \phi(a, x)^\top]^{-1} \mathbb{E}_n^{(1)} [\phi(a, x) \partial_h m(a, x; \theta, h, g)],$$

and $\partial_h m(\cdot)$ denotes the Fréchet derivative of $m(\cdot)$ with respect to $h(\cdot)$. The aim is to obtain an expression for the limit, \mathcal{A}_h , of $\hat{\mathcal{A}}_h$ as $n, k_\phi \rightarrow \infty$. We will then conjecture \mathcal{A}_h to be the non-parametric correction term for $h(\cdot)$.

To this end, let us keep the dimension k_ϕ fixed for now and define

$$\hat{\vartheta} := \mathbb{E}_n^{(1)} [(\beta \phi(a', x') - \phi(a, x)) \phi(a, x)^\top]^{-1} \mathbb{E}_n^{(1)} [\phi(a, x) \partial_h m(a, x; \theta, h, g)].$$

Note that $\hat{\lambda}_h(a, x) = \phi(a, x)^\top \hat{\vartheta}$. Now, in the limit as $n \rightarrow \infty$, we can expect $\hat{\vartheta} - \vartheta \rightarrow 0$, where

$$\vartheta := \mathbb{E} [(\beta \phi(a', x') - \phi(a, x)) \phi(a, x)^\top]^{-1} \mathbb{E} [\phi(a, x) \partial_h m(a, x; \theta, h, g)].$$

Since $\mathbb{E}[\cdot]$ is a stationary distribution, $\mathbb{E} [\beta \phi(a', x') \phi(a, x)^\top] = \mathbb{E} [\beta \phi(a, x) \phi(a^-, x^-)^\top]$, where (a^-, x^-) denotes the one-step backward quantities corresponding to (a, x) . In view of this, a bit of rearrangement of the previous display equation gives us

$$\mathbb{E} [\phi(a, x) \{-\partial_h m(a, x; \theta, h, g) + \beta \phi(a^-, x^-)^\top \vartheta - \phi(a, x)^\top \vartheta\}] = 0. \quad (3.16)$$

Define $\lambda_h^*(a, x) := \phi(a, x)^\top \vartheta$, noting also that this is the limit of $\hat{\lambda}_h(a, x)$ as $n \rightarrow \infty$. Given (3.16), we then have

$$\mathbb{E} [\phi(a, x) \{-\partial_h m(a, x; \theta, h, g) + \beta \lambda_h^*(a^-, x^-) - \lambda_h^*(a, x)\}] = 0.$$

The above equation shares a high degree of similarity with (3.1). Indeed, backtracking the analysis leading to (3.1), we see that $\lambda_h^*(a, x)$ can be interpreted as the fixed point of the projected ‘backward’ dynamic programming operator $P_\phi \Gamma_h^\dagger[\cdot]$, where⁵

$$\Gamma_h^\dagger[f](a, x) := -\partial_h m(a, x; \theta, h, g) + \beta \mathbb{E} [f(a^-, x^-) | a, x].$$

While we have supposed the dimension of $\phi(\cdot)$ to be fixed so far, as $k_\phi \rightarrow \infty$, we can expect $\lambda_h^*(a, x) \rightarrow \lambda_h(a, x)$, where the latter is the fixed point of $\Gamma_h^\dagger[\cdot]$ itself. From the above discussion, we can conjecture that the limit of $\hat{\mathcal{A}}_h$ is given by

$$\mathcal{A}_h = \lambda_h(a, x) \{z(a, x) + \beta h(a', x') - h(a, x)\},$$

where we have also replaced $\phi(a, x)^\top \omega$ in (3.15) with its limit $h(a, x)$. This is our conjecture for the adjustment term corresponding to $h(\cdot)$. A similar analysis also applies to the adjustment

⁵In other words, $\lambda_h^*(a_{it}, x_{it}) = -\sum_{j=0}^{\infty} \beta^j \partial_h m(a_{i(t-j)}, x_{i(t-j)}; \theta, h, g)$, i.e it is essentially like a ‘backward’ value function.

term for $g(\cdot)$, which we conjecture to be of the form

$$\mathcal{A}_g = \lambda_g(a, x) \{e(a', x'; \eta) + \beta g(a', x') - g(a, x)\},$$

where $\lambda_g(a, x)$ is the fixed point of the operator $\Gamma_g^\dagger[\cdot]$, defined as

$$\Gamma_g^\dagger[f](a, x) := -\partial_g m(a, x; \theta, h, g) + \beta \mathbb{E}[f(a^-, x^-) | a, x].$$

Taken together, we conjecture that the locally robust moment is given by

$$\begin{aligned} \zeta(a, x; \theta, h, g) &:= m(a, x; \theta, h, g) - \lambda_h(a, x) \{z(a, x) + \beta h(a', x') - h(a, x)\} \\ &\quad - \lambda_g(a, x) \{e(a', x'; \eta) + \beta g(a', x') - g(a, x)\}. \end{aligned} \quad (3.17)$$

The above analysis is heuristic. We now verify that the moment in (3.17) is indeed locally robust. A necessary condition for this is that $\partial_h \mathbb{E}[\zeta(a, x; \theta, h, g)] = 0$ and $\partial_g \mathbb{E}[\zeta(a, x; \theta, h, g)] = 0$, where the derivatives are Gâteaux derivatives with respect to $h(\cdot)$ and $g(\cdot)$ respectively (see Chernozhukov et al. (2018)). To verify these, observe that for any square integrable γ ,

$$\begin{aligned} \partial_\tau \mathbb{E}[\zeta(a, x; \theta, h + \tau \gamma, g)] &= \mathbb{E}[\partial_h m(a, x; \theta, h, g) \gamma(a, x)] - \mathbb{E}[\beta \lambda_h(a, x) \gamma(a', x')] \\ &\quad + \mathbb{E}[\lambda_h(a, x) \gamma(a, x)]. \end{aligned} \quad (3.18)$$

Since $\lambda_h(\cdot)$ is the fixed point of $\Gamma_h^\dagger[\cdot]$, we can expand the third term in (3.18) as

$$\begin{aligned} \mathbb{E}[\lambda_h(a, x) \gamma(a, x)] &= \mathbb{E}[\{-\partial_h m(a, x; \theta, h, g) + \beta \lambda_h(a^-, x^-)\} \gamma(a, x)] \\ &= -\mathbb{E}[\partial_h m(a, x; \theta, h, g) \gamma(a, x)] + \mathbb{E}[\beta \lambda_h(a, x) \gamma(a', x')], \end{aligned}$$

where the second equality uses the fact that $\mathbb{E}[\cdot]$ is a stationary distribution. We thus conclude $\partial_\tau \mathbb{E}[\zeta(a, x; \theta, h + \tau \gamma, g)] = 0$ for all γ , or $\partial_h \mathbb{E}[\zeta(a, x; \theta, h, g)] = 0$, as required. In fact, by a similar argument to the above, we can also show the stronger statement that $\partial_h \mathbb{E}[\zeta(a, x; \theta, h, g)] = 0$ and $\partial_g \mathbb{E}[\zeta(a, x; \theta, h, g)] = 0$ in a Fréchet sense. Additionally, the Fréchet second derivatives $\partial_h^2 \mathbb{E}[\zeta(a, x; \theta, h, g)]$ and $\partial_g^2 \mathbb{E}[\zeta(a, x; \theta, h, g)]$ also exist, and are uniformly bounded over all θ lying in a compact set.

The locally robust moment (3.17) is infeasible since $\lambda_g(\cdot)$, $\lambda_h(\cdot)$, $h(\cdot)$ and $g(\cdot)$ are unknown. However, in practice we can simply use the estimator from (3.14). Note that the moment function from the latter can be rewritten as

$$\begin{aligned} \zeta_n(a, x; \theta, \mathbf{v}) &= m(a, x; \theta, \mathbf{v}) - \hat{\lambda}_h(a, x) \{z(a, x) + \beta \phi(a', x')^\top \omega - \phi(a, x)^\top \omega\} \\ &\quad - \hat{\lambda}_g(a, x) \{e(a', x'; \hat{\eta}) + \beta r(a', x')^\top \xi - r(a, x)^\top \xi\}. \end{aligned}$$

There is no loss of first order efficiency in replacing $\zeta(a, x; \theta, h, g)$ with $\zeta_n(a, x; \theta, \mathbf{v})$. This is because, by a similar analysis as for Theorems 1, 2, it can be shown that

$$\begin{aligned} \|\hat{\lambda}_h(a, x) - \lambda_h(a, x)\|_2 &= O_p\left(\frac{k_\phi}{\sqrt{n}} + k_\phi^{-\alpha}\right) = o_p(n^{-1/4}), \text{ and} \\ \|\hat{\lambda}_g(a, x) - \lambda_g(a, x)\|_2 &= O_p\left(\frac{k_r}{\sqrt{n}} + k_r^{-\alpha}\right) = o_p(n^{-1/4}), \end{aligned}$$

for suitable choices of k_r and k_ϕ . Note also that $\phi(a, x)^\top \omega$ and $r(a, x)^\top \xi$ are L_2 consistent for $h(a, x)$ and $g(a, x)$, respectively, at faster than $n^{-1/4}$ rates. Following the analysis of Chernozhukov et al. (2018), these facts imply that the estimator based on (3.13) has the same limiting distribution as the one based on (3.17). In particular, it achieves parametric rates of convergence. We state the regularity conditions and the theorem below (for the remainder of this section we allow θ^* to be vector valued):

Assumption 3. (i) $\theta^* \in \Theta$, where Θ is a compact set.

(ii) $\partial_g m(a, x; \theta, h, g)$ and $\partial_h m(a, x; \theta, h, g)$ are uniformly bounded for all (a, x, θ) .

(iii) Let $G := \mathbb{E} [\partial_\theta \zeta(a, x; \theta^*, h, g)]$. Then G is invertible.

Theorem 3. Suppose that Assumptions 1 - 3 hold. Then the estimator, $\hat{\theta}$ of θ^* , based on (3.13) is \sqrt{n} consistent, and satisfies

$$\sqrt{n}(\hat{\theta} - \theta^*) \implies N(0, V),$$

where $V = (G^\top \Omega^{-1} G)^{-1}$, with $\Omega := \mathbb{E} [\zeta(a, x; \theta^*, h, g) \zeta(a, x; \theta^*, h, g)^\top]$.

The proof of the above theorem follows by verifying the regularity conditions of Chernozhukov et al. (2018, Theorem 16). Since these are more or less straightforward to verify given our previous results, we omit the details.

For inference on $\hat{\theta}$, the covariance matrix V can be estimated as

$$\hat{V} = \left(\hat{G}^\top \hat{\Omega}^{-1} \hat{G} \right)^{-1},$$

where

$$\begin{aligned} \hat{G} &= \frac{1}{n(T-1)} \sum_{i=1}^n \sum_{t=1}^{T-1} \frac{\partial \zeta_n(a_{it}, x_{it}; \hat{\theta}, \hat{\omega}, \hat{\xi})}{\partial \theta^\top}, \quad \text{and} \\ \hat{\Omega} &= \frac{1}{n(T-1)} \sum_{i=1}^n \sum_{t=1}^{T-1} \zeta_n(a_{it}, x_{it}; \hat{\theta}, \hat{\omega}, \hat{\xi}) \zeta_n(a_{it}, x_{it}; \hat{\theta}, \hat{\omega}, \hat{\xi})^\top. \end{aligned}$$

Chernozhukov et al. (2018) provide conditions under which \hat{V} is consistent for V ; these are straightforward to verify in our context. Alternatively, one could employ the bootstrap.

4. EXTENSIONS

4.1. Stochastic Gradient descent. Computation of $\hat{\omega}$ and $\hat{\xi}$ involves inverting a $(k \times k)$ -dimensional matrix. Once k becomes very large, matrix inversion does start to become more demanding. In such cases stochastic gradient descent is a computationally cheap alternative. In particular, we can estimate ω^* in (3.3) using stochastic gradient updates of the form

$$\hat{\omega}^{new} \leftarrow \hat{\omega}^{old} + \alpha_\omega \left(z_{it} + \beta \phi_{it+1}^\top \hat{\omega}^{old} - \phi_{it}^\top \hat{\omega}^{old} \right) \phi_{it}, \quad (4.1)$$

where each observation $(z_{it}, \phi_{it}, \phi_{it+1})$ is drawn at random from \mathbb{P}_n i.e, with replacement from the set of all the sample observations. Here α_ω is the learning rate for stochastic gradient descent. In a similar vein we can estimate ξ^* using gradient updates of the form

$$\hat{\xi}^{new} \leftarrow \hat{\xi}^{old} + \alpha_\xi \left(\beta e_{it+1}(\hat{\eta}) + \beta r_{it+1}^\top \hat{\xi}^{old} - r_{it}^\top \hat{\xi}^{old} \right) r_{it}, \quad (4.2)$$

Algorithm 1 TD learning algorithm for CCP estimation

Initialize all parameters to arbitrary values

Repeat:

Choose $(x_{it}, a_{it}, x_{it+1}, a_{it+1})$ at random, with replacement, from sample data

Calculate the values of $(\phi_{it}, z_{it}, r_{it}, \phi_{it+1}, r_{it+1}, e_{it+1}(\hat{\eta}))$

$$\hat{\omega} \leftarrow \hat{\omega} + \alpha_{\omega} \left(z_{it} + \beta \phi_{it+1}^{\top} \hat{\omega} - \phi_{it}^{\top} \hat{\omega} \right) \phi_{it}$$

$$\hat{\xi} \leftarrow \hat{\xi} + \alpha_{\xi} \left(\beta e_{it+1}(\hat{\eta}) + \beta r_{it+1}^{\top} \hat{\xi} - r_{it}^{\top} \hat{\xi} \right) r_{it}$$

Until: Convergence criteria for $(\hat{\omega}, \hat{\xi})$ are reached

where α_{ξ} is the learning rate for ξ , and $e_{it+1}(\hat{\eta}) := \gamma - \ln \hat{\eta}(a_{it+1}, x_{it+1})$. Estimation of (ω^*, ξ^*) using the gradient updates (4.1) and (4.2) is termed TD learning in the Reinforcement Learning literature. Pseudo-code for our TD learning algorithm is provided in Algorithm 1.

We shall require the following assumption on the learning rates:

Assumption 4. *The learning rates satisfy $\sum_l \alpha_{\omega}^{(l)2} \rightarrow 0$, $\sum_l \alpha_{\xi}^{(l)2} \rightarrow 0$ and $\sum_l \alpha_{\omega}^{(l)} \rightarrow \infty$, $\sum_l \alpha_{\xi}^{(l)} \rightarrow \infty$ as the number of steps in the algorithm goes to infinity, where $\alpha_{\omega}^{(l)}, \alpha_{\xi}^{(l)}$ denote the learning rates after l steps/updates of the algorithm.*

Assumption 4 is a standard condition on learning rates for stochastic gradient descent algorithms. We can now prove the following theorem on convergence:

Theorem 4. *Suppose that Assumptions 1, 2 and 4 hold. Then, with probability approaching one, the sequence of updates ω_l and ξ_l converge to $\hat{\omega}, \hat{\xi}$ as $l \rightarrow \infty$.*

The TD learning algorithm can also be parallelized by running multiple stochastic gradient threads in parallel and using Hogwild!-style asynchronous updates (Niu et al., 2011). Each thread runs parallel instances of the same code with a delayed time start, and independently and asynchronously updates a global parameter that returns ω . This speeds up computation by the order of magnitude of the number of parallel threads.

4.2. Nonlinear utility functions. So far we have focused on the case where the utility function is linear in parameters θ^* . This simplifies the computation considerably. However, in practice it may be useful to specify the observed utility component to be nonlinear in the parameters θ^* . Denote this by $z(a, x; \theta^*)$. In such cases, we may estimate θ^* as the maximizer of the quasi-likelihood criterion

$$Q(\theta) = \sum_{i=1}^n \sum_{t=1}^T \log \frac{\exp \{h(a_{it}, x_{it}; \theta) + g(a_{it}, x_{it})\}}{\sum_a \exp \{h(a, x_{it}; \theta) + g(a, x_{it})\}},$$

where, for each θ , $h(\cdot; \theta)$ and $g(\cdot)$ solve the following recursive expressions:

$$h(a, x; \theta) = z(a, x; \theta) + \beta \mathbb{E} [h(a', x'; \theta) | a, x], \text{ and}$$

$$g(a, x) = \beta \mathbb{E} [e(a', x') + g(a', x') | a, x].$$

We can use our TD estimation procedure to obtain a functional approximation $\hat{h}(a, x; \theta)$ for $h(a, x; \theta)$, conditional on each different value of θ . As argued earlier, this step only requires a low-dimensional matrix inversion and can be computed extremely fast. As long as Assumption 1 holds uniformly over all θ , we can also prove that $\hat{h}(a, x; \theta)$ is uniformly consistent for $h(a, x; \theta)$ at the same rates as before i.e.

$$\sup_{\theta \in \Theta} \|h(a, x; \theta) - \hat{h}(a, x; \theta)\|_2 = O_p \left(\frac{k_r}{\sqrt{n}} + k_r^{-\alpha} \right).$$

We can therefore plug in the values of $\hat{h}(\cdot; \theta)$ and $\hat{g}(\cdot)$ to estimate θ^* as

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \hat{Q}(\theta); \quad \hat{Q}(\theta) := \sum_{i=1}^n \sum_{t=1}^T \log \frac{\exp \{ \hat{h}(a_{it}, x_{it}; \theta) + \hat{g}(a_{it}, x_{it}) \}}{\sum_a \exp \{ \hat{h}(a, x_{it}; \theta) + \hat{g}(a, x_{it}) \}}.$$

The main computational difficulty lies in maximizing $\hat{Q}(\theta)$ with respect to θ . Since the derivatives cannot be computed in straightforward manner, we recommend a derivative-free optimization procedure such as Nelder-Mead to solve for $\hat{\theta}$. The large sample properties of $\hat{\theta}$ are more involved and we do not attempt to derive them here.

4.3. Recursive estimation. A drawback of the estimation strategy so far is that it still requires some initial estimates of the choice probabilities to obtain the values of $e(a', x')$. This is so even as we do eliminate entirely the need for any initial probability values when estimating $h(\cdot)$, as well as the need to estimate f_X . In this section we show how the estimate for η can also be dispensed with, at the expense of a bit more computation. The key insight we exploit is the fact that at the true value θ^* of θ , we will have

$$\eta(a, x) = \frac{\exp \{ h(a, x) \theta^* + g(a, x) \}}{\sum_a \exp \{ h(a, x) \theta^* + g(a, x) \}}.$$

Thus, if we have a consistent estimator for θ^* , we can use this to obtain an estimate for $\eta(a, x)$. This suggests a recursive procedure for estimating $\eta(\cdot)$ and θ simultaneously.

Note that, even with this recursive procedure, the estimates $\hat{\omega}$ can be obtained directly from (3.3). We do not require any estimate of $\eta(a, x)$ for this. Let $\hat{h}(a, x)$ denote the estimate of $h(a, x)$ that we obtained in the previous sections. We start the recursive procedure by initializing ξ and θ to arbitrary values. Additionally, we also initialize $\eta(a, x)$ by $\hat{\eta}_{(1)}(a, x)$, where the latter is some preliminary estimate of the choice probabilities. Let $\hat{\xi}_{(k)}$ and $\hat{\theta}_{(k)}$ denote the parameter estimates, at the k -th iteration of the procedure. Similarly, let $\hat{\eta}_{(k)}(a, x)$ and $\hat{e}_{(k)}(a, x)$ denote the estimates of $\eta(\cdot)$ and $e(\cdot)$ after k iterations of the procedure. These quantities are then updated as follows: We first update $\hat{\eta}(\cdot)$ as

$$\hat{\eta}_{(k+1)}(a, x) = \frac{\exp \{ \hat{h}(a, x) \hat{\theta}_{(k)} + r(a, x)^\top \hat{\xi}_{(k)} \}}{\sum_{\hat{a}} \exp \{ \hat{h}(\hat{a}, x) \hat{\theta}_{(k)} + r(\hat{a}, x)^\top \hat{\xi}_{(k)} \}}. \quad (4.3)$$

This enables us to obtain a new estimate of $e(a, x)$,

$$\hat{e}_{(k+1)}(a, x) := \gamma - \ln \hat{\eta}_{(k+1)}(a, x). \quad (4.4)$$

Following this, $\hat{\xi}$ can be updated as

$$\hat{\xi}_{(k+1)} = \mathbb{E}_n [r(a, x) (r(a, x) - \beta r(a', x'))^\top]^{-1} \mathbb{E}_n [\beta r(a, x) \hat{e}_{(k+1)}(a', x')]. \quad (4.5)$$

Finally, $\hat{\theta}$ can be updated as

$$\hat{\theta}_{(k+1)} = \arg \max_{\theta} \sum_{i=1}^n \sum_{t=1}^{T-1} \log \frac{\exp \{ \hat{h}(a_{it}, x_{it}) \theta + r(a_{it}, x_{it})^\top \hat{\xi}_{(k+1)} \}}{\sum_a \exp \{ \hat{h}(a, x_{it}) \theta + r(a, x_{it})^\top \hat{\xi}_{(k+1)} \}}. \quad (4.6)$$

The above update does not employ the locally robust correction to obtain $\hat{\theta}$. This can be easily rectified using (3.14); we refer to the previous section for the details. We iterate between steps (4.3) - (4.6) until the parameters converge.

Our recursive procedure is very similar to, and influenced by, the NPL algorithm of Aguirregabiria and Mira (2002). Using Monte Carlo simulations, the authors show that the recursive procedure enjoys smaller finite sample bias and variance. This was subsequently proven using higher order expansions by Kasahara and Shimotsu (2008). We similarly expect our recursive procedure to have better finite sample properties. Furthermore, as with the NPL algorithm, it can be shown that in the case of discrete states, the recursive procedure converges to the Maximum Likelihood estimate obtained using the Rust (1987) NFXP algorithm as the number of iterations increases to infinity.

5. INCORPORATING PERMANENT UNOBSERVED HETEROGENEITY

In this section, we show how we can model permanent unobserved heterogeneity by pairing the techniques from Section 3 with the Expectation-Maximization (EM) algorithm. The use of the EM algorithm in CCP estimation under unobserved heterogeneity was first advocated by Arcidiacono and Miller (2011), and we employ the same approach.

Suppose that in addition to the observed state x , and the choice specific shock e , individuals also base their choice decisions on a random state variable s which is known to the individual, but unobserved to the econometrician. As is common in the literature, we assume a finite set of unobserved states indexed by $\{1, 2, \dots, k, \dots, K\}$. The number of states is also assumed to be known a priori. Let π_k denote the population probability $P(s = k)$. The value of s for an individual is assumed to be permanent and not change with time. We therefore treat the possible realizations of s as indices in the value function approximation where, conditional on realizations (a_{it}, x_{it}) , i.e. each individual has K potential value functions.

We shall also make two further simplifications to ease the exposition: First, we suppose that all the states are really discrete. The same procedure can also be applied to continuous states, but it is not efficient in this case. To regain efficiency, we have to employ local robustness corrections as in Section 2. Second, we make a random effects assumption that the unobserved state variable is independent of the observed states $\{x_1, \dots, x_T\}$. This enables us to avoid specifying a likelihood function for the observed states.

To simplify notation, similar to before, let $h_{itk} := h_k(a_{it}, x_{it})$ and $g_{itk} := g_k(a_{it}, x_{it})$. Suppose that the population probabilities π_k are known. Then one can estimate the structural parameters

θ by maximizing the integrated likelihood⁶

$$Q(\theta) = \sum_{i=1}^N \log \left[\sum_{k=1}^K \pi_k \prod_{t=1}^T \frac{\exp \{h_{itk}\theta + g_{itk}\}}{\sum_a \exp \{h_k(a, x_{it})\theta + g_k(a, x_{it})\}} \right]. \quad (5.1)$$

Since the value functions now depend on s , for our functional approximations, we choose $\phi(a, x, k)$ and $r(a, x, k)$ as a set of basis functions over the domain (a, x, k) so that we may approximately parameterize

$$h_k(a, x) = \phi_k(a, x)^\top \omega^*; \quad g_k(a, x) = r_k^\top(a, x) \xi^*.$$

As before, we have chosen to make $h()$ uni-dimensional to simplify the notation. The extension to multiple dimensions is straightforward as one simply treats each dimension separately.

For the case of known π_k , we can modify our earlier procedure as follows (the updates for ξ are similar): Similar to Section 3, ω^* is identified as

$$\omega^* = \bar{\mathbb{E}} [\phi_k(a, x) (\phi_k(a, x) - \beta \phi_k(a', x'))^\top]^{-1} \bar{\mathbb{E}} [\phi_k(a, x) z_k(a, x)],$$

where $\bar{\mathbb{E}}[\cdot]$ differs from $\mathbb{E}[\cdot]$ in also taking the expectation over the distribution of the unobserved state s . In particular, observe that

$$\bar{\mathbb{E}} [\phi_k(a, x) z_k(a, x)] = \mathbb{E} \left[\sum_k P(s = k | \mathbf{x}, \mathbf{a}) \phi_k(a, x) z_k(a, x) \right],$$

where

$$P(s = k | \mathbf{x}, \mathbf{a}) := Pr(s = k | x_1, \dots, x_T, a_1, \dots, a_T)$$

denotes the probability that the unobserved state is k conditional on the set of all the actions and observed states for an individual. Note that the last equation in the above expression follows by the law of iterated expectations. In a similar vein, we also have

$$\begin{aligned} & \bar{\mathbb{E}} [\phi_k(a, x) (\phi_k(a, x) - \beta \phi_k(a', x'))^\top] \\ &= \mathbb{E} \left[\sum_k P(s = k | \mathbf{x}, \mathbf{a}) \phi_k(a, x) (\phi_k(a, x) - \beta \phi_k(a', x'))^\top \right]. \end{aligned}$$

Denote by $p_{ik} = P(s = k | \mathbf{x}_i, \mathbf{a}_i)$ the probability of being in state k conditional on the realized set of all actions \mathbf{a}_i and observed states \mathbf{x}_i for individual i .

To further simplify notation, denote $z_{itk} := z_k(a_{it}, x_{it})$, $\phi_{itk} := \phi_k(a_{it}, x_{it})$, $e_{itk} := e_k(a_{it}, x_{it})$, and $r_{itk} := r_k(a_{it}, x_{it})$. Then replacing the expectation $\mathbb{E}[\cdot]$ in the previously displayed equations with the sample expectation $\mathbb{E}_n[\cdot]$, we obtain the estimate

$$\hat{\omega} = \left[\sum_{i=1}^n \sum_{t=1}^{T-1} \sum_k p_{ik} \phi_{itk} (\phi_{itk} - \beta \phi_{it+1k})^\top \right]^{-1} \sum_{i=1}^n \sum_{t=1}^{T-1} \sum_k p_{ik} \phi_{itk} z_{itk} \quad (5.2)$$

A similar expression also holds for updates to ξ :

$$\hat{\xi} = \left[\sum_{i=1}^n \sum_{t=1}^{T-1} \sum_k p_{ik} r_{itk} (r_{itk} - \beta r_{it+1k})^\top \right]^{-1} \sum_{i=1}^n \sum_{t=1}^{T-1} \sum_k \beta p_{ik} r_{itk} e_{it+1k}, \quad (5.3)$$

⁶This is in fact a conditional likelihood (i.e conditional on $\{x_{i1}, \dots, x_{iT}\}$) since we assumed s is independent of the observed states.

where \dot{e}_{it+1k} is the current estimate of e_{it+1k} .

Estimation of ω^*, ξ^* and θ^* using equations (5.1), (5.2) and (5.3) requires knowledge of the unknown quantities π_k and p_{ik} along with \dot{e}_{it+1k} . Furthermore, even if π_k were known, maximizing the integrated likelihood function (5.1) is computationally very expensive. The EM algorithm solves both issues and provides a computationally cheap alternative to maximizing (5.1). We modify the EM algorithm slightly to additionally include updates to the estimate \dot{e}_{it+1k} of e_{it+1k} , drawing on ideas from Section 4.3. To describe the procedure, let

$$l_{itk}(\theta, \omega, \xi) \equiv \frac{\exp\{(\phi_{itk}^\top \omega)\theta + (r_{itk}^\top \xi)\}}{\sum_a \exp\{(\phi_k(a, x_{it})^\top \omega)\theta + r_k(a, x_{it})^\top \xi\}}.$$

Denote by $\hat{\pi}_k$ and \hat{p}_{ik} the estimates for π_k and p_{ik} . The algorithm consists of two steps: the M-step and the E-step. We first describe the M-step. Here, we update the estimates for ω^*, ξ^* and θ^* based on the current estimates for π_k, p_{ik} and e_{it+1k} . To this end, first note that we can update $\hat{\omega}$ and $\hat{\xi}$ using (5.2) and (5.3). From these we can in-turn update $\hat{\theta}$ as

$$\hat{\theta} = \arg \max_{\theta} \left[\sum_k p_{ik} \ln l_{itk}(\theta, \hat{\omega}, \hat{\xi}) \right]. \quad (5.4)$$

Next, given $\hat{\theta}, \hat{\omega}$ and $\hat{\xi}$, we update $\hat{\pi}_k, \hat{p}_{ik}$ and \dot{e}_{it+1k} for all i, k . This is the E-step of the EM algorithm. This step consists of three parts. In the first part, we use the current $\hat{\theta}, \hat{\omega}, \hat{\xi}$ and $\hat{\pi}_k$ to update \hat{p}_{ik} for each i, k using Bayes' rule:

$$\hat{p}_{ik} \leftarrow \frac{\hat{\pi}_k \prod_{t=1}^T l_{itk}(\hat{\theta}, \hat{\omega}, \hat{\xi})}{\sum_{\tilde{k}} \hat{\pi}_{\tilde{k}} \prod_{t=1}^T l_{it\tilde{k}}(\hat{\theta}, \hat{\omega}, \hat{\xi})}. \quad (5.5)$$

In the second part, we update $\hat{\pi}_k$, for each k , as

$$\hat{\pi}_k \leftarrow \frac{1}{N} \sum_{i=1}^N \hat{p}_{ik}. \quad (5.6)$$

Finally, we also update \dot{e}_{it+1k} for all i, t, k as

$$\dot{e}_{it+1k} \leftarrow \gamma - \ln l_{it+1k}(\hat{\theta}, \hat{\omega}, \hat{\xi}), \quad (5.7)$$

recalling that $\mathcal{G}(\cdot)$ is the function mapping the conditional probabilities to the error term.

The computational requirements for the EM algorithm are higher due to the iteration between the expectation and maximization steps. However the maximization step is still very fast as we can estimate all the parameters ω^* and ξ^* through a low-dimensional matrix inversion, while obtaining $\hat{\theta}$ is just a convex optimization problem. Furthermore, as in Section 3, one can use stochastic gradient descent as an alternative to matrix inversion.

6. ESTIMATION OF DYNAMIC DISCRETE GAMES

So far we have considered applications of our algorithm to single agent models, where we have argued that there are substantial computational and statistical gains from using our procedure. These gains are magnified when extended to estimation of dynamic discrete games.

Our setup is based on Aguirregabiria and Mira (2010). We assume a single Markov-Perfect-Equilibrium setup where multiple players $i = 1, 2, \dots, N$ play against each other in M different

markets. We observe the state of play for T time-periods, where $T \ll N$. Utility of the players in any time period is affected by the actions of all the others, and a set of states x that are observed by all players. The per period utility is denoted by $z_i(a_i, a_{-i}, x)^\top \theta^*$ for each player i , for some finite dimensional parameter θ^* , where a_i denotes player i 's action and a_{-i} denotes the actions of all other players. Evolution of the states in the next period is determined by the transition probability $f_X(x'|a, x)$ where $\mathbf{a} := (a_1, \dots, a_N)$ denotes the actions of all the players. We denote by x_{tm} the state at market m in time period t , and by a_{itm} the action of player i at time t in market m . We also let $P_i(a_i|x_t)$ denote the choice probability of player i taking action a_i when the state is x_t .

As in the single agent case, the parameters θ^* can be obtained as solutions to the pseudo-likelihood function:

$$Q(\theta) = \sum_{m=1}^M \sum_{i=1}^N \sum_{t=1}^T \log \frac{\exp \{h(a_{itm}, x_{tm})^\top \theta + g(a_{itm}, x_{tm})\}}{\sum_a \exp \{h(a, x_{tm})^\top \theta + g(a, x_{tm})\}}, \quad (6.1)$$

where $h(\cdot)$ and $g(\cdot)$ are defined very similarly to $h(\cdot)$ and $g(\cdot)$ in the single agent case, the complication being that actions of other players need to be partialled out:

$$\begin{aligned} h(a_i, x) &= \sum_{a_{-i}} \left(\prod_{j \neq i} P_j(a_j|x) \right) \left[z(a_i, a_{-i}, x) + \beta \sum_{x'} f_X(x'|a_i, a_{-i}, x) \sum_{a'} P_i(a'|x') h(a', x') \right] \\ g(a_i, x) &= \beta \sum_{a_{-i}} \left(\prod_{j \neq i} P_j(a_j|x) \right) \left[\sum_{x'} f_X(x'|a_i, a_{-i}, x) \sum_{a'} P_i(a'|x') \{e(a', x') + g(a', x')\} \right]. \end{aligned}$$

Converting the above to expectations gives us

$$\begin{aligned} h(a_i, x) &= \mathbb{E}[z(a_i, a_{-i}, x)|a_i, x] + \beta \mathbb{E}[h(a', x')|a_i, x], \\ g(a_i, x) &= \mathbb{E}[e(a', x') + \beta g(a', x')|a_i, x]. \end{aligned} \quad (6.2)$$

In contrast to (2.1) in the single agent case, the expectation now averages over the actions of the other players as well.

Previous literature estimates θ^* using a two-step procedure: In the first step, the conditional choice probabilities $P_i(a_i|x_t)$ are calculated non-parametrically. These, along with estimates of $f_X(\cdot)$ are then used to recursively solve for $h(\cdot)$ and $g(\cdot)$ using equation (6.2). This step requires integrating over the actions of all the other players. Finally, given the estimated values of $h(\cdot)$ and $g(\cdot)$, the parameter θ is estimated through either pseudo-likelihood (Aguirregabiria and Mira, 2007) or minimum distance estimation (Pesendorfer and Schmidt-Dengler, 2008). Both these approaches have been proposed for discrete states. For continuous states, Bajari et al. (2007) have proposed an alternative method to solve (6.2) by forward Monte Carlo simulation. Though computationally cheaper than discretization (which could give rise to a very high dimension of states), forward simulation is still cumbersome with many continuous states and players.

By contrast, our algorithm is a straightforward extension of the ones we suggested in earlier sections for single agent models. Let $\hat{\eta}(a, x)$ denote some non-parametric estimate of the choice probabilities. Here, and in what follows, we shall use a to denote the action of any particular

individual player. We then (approximately) parameterize $h(\cdot)$ and $g(\cdot)$ as

$$h(a, x) : \approx \phi(a, x)^\top \omega^*; \quad g(a, x) : \approx r(a, x)^\top \xi^*,$$

where, as before, $\phi(a, x)$ and $r(a, x)$ are comprised of a set of basis functions over the domain of (a, x) .

We now simply proceed to estimate the value weights ω^* and ξ^* exactly as in Section 3:

$$\begin{aligned} \hat{\omega} &= \mathbb{E}_n [\phi(a, x) (\phi(a, x) - \beta \phi(a', x'))^\top]^{-1} \mathbb{E}_n [\phi(a, x) z(a, x)], \\ \hat{\xi} &= \mathbb{E}_n [r(a, x) (r(a, x) - \beta r(a', x'))^\top]^{-1} \mathbb{E}_n [\beta r(a, x) e(a', x'; \hat{\eta})] \end{aligned} \quad (6.3)$$

The functions $h(\cdot)$ and $g(\cdot)$ are common for all the players. This is without loss of generality, however, as one could always decompose θ, ω, ξ into player-specific components, e.g. $\theta = (\theta_1, \dots, \theta_N)$. Note that one could also apply the procedure separately for each player, in which case we would replace $h(\cdot)$ and $g(\cdot)$ in (6.1) with $h_m(\cdot)$ and $g_m(\cdot)$ to reflect the fact that these quantities are now player specific.

Remarkably, the estimation strategy in (6.3) does not require partialling out the other players' actions, leading to a tremendous reduction of computation. Indeed, the procedure automatically takes expectations over the actions of the other players using the empirical distribution. To see this in the discrete case, note that \bar{z}_{itm} is an unbiased estimator of the expectation $\sum_{a_{-i}} \prod_{j \neq i} P_j(a_j | x) z(a_i, a_{-i}, x)$. This intuition also goes through with continuous states since we use a functional approximation, which provides an automatic regularization for calculating the above expectation 'internally' as long as the dimension of $\phi(\cdot)$ and $r(\cdot)$ is sufficiently small relative to the sample size.

By the same reasoning as in Section 3.1, it is possible to show that with discrete states, $h(\cdot)$ and $g(\cdot)$ are numerically identical to the estimates obtained by plugging in cell estimates $\hat{P}_j(\cdot | x)$ and $\hat{f}_X(\cdot)$ in (6.2). This implies the partial likelihood with plug-in estimates for $h(\cdot)$ and $g(\cdot)$ is not efficient even with discrete states, as discussed by Aguirregabiria and Mira (2007). However the values of $h(\cdot)$ and $g(\cdot)$ can be plugged into other, more efficient objectives, such as the minimum distance estimator of Pesendorfer and Schmidt-Dengler (2008). With continuous states, one would need to employ locally robust corrections to recover parametric rates of convergence for θ . To this end, we can use the fact that the minimum distance estimator can be characterized by a moment criterion. Combining this with the moments implied by (6.3) for ω and ξ , it is easy to see how the construction of Section 3.3 can be extended to dynamic games. One could also use a recursive version of our algorithm as in Section 4.3. This is equivalent to full information MLE under some additional conditions (see Kasahara and Shimotsu, 2012).

Finally, we note that it is also straightforward to incorporate the other extensions from Section 4 to the setup of dynamic games..

7. SIMULATIONS

We run Monte Carlo Simulations to test our estimation method, starting with the simplest version of our algorithm described in Section 3. Our simulations are based on a modified version

of the Rust (1987) engine replacement problem. We start by describing the setup in Section 7.1, before moving to the simulation results in Section 7.2.

7.1. Bus Engine Replacement Problem. Consider the following version of the Rust (1987) bus engine replacement problem which is adapted from Arcidiacono and Miller (2011). Each period $t = 1, \dots, T; T < \infty$, Harold Zurcher decides whether to replace the engine of a bus ($a_t = 0$), or keep it ($a_t = 1$). Denote his action by $j \in \{0, 1\}$. Each bus is characterized by a permanent type $s \in \{1, 2\}$, and the mileage accumulated since the last engine replacement $x_t \in \{1, 2, \dots\}$. Harold Zurcher observes both s and x_t . As in Section 3, we start by also treating both s and x_t as observed to the econometrician.

Mileage increases by one unit if the engine is kept in period t and is set to zero if the engine is replaced. The current period payoff for keeping the engine is given by $\theta_0 + \theta_1 x_t + \theta_2 s + e_{1t}$, where $\theta^* \equiv \{\theta_0, \theta_1, \theta_2\}$ are the structural parameters of interest, and e_{jt} is a choice-specific transitory shock that follows a Type 1 Extreme Value distribution. As in Arcidiacono and Miller (2011), we normalize the current period payoff of replacing the engine to e_{0t} .

When deciding whether to keep or replace the engine, Harold Zurcher solves a DDC problem and sequentially maximizes the following discounted sum of payoffs:

$$E \left[\sum_{t=1}^T \beta^t \{a_t(\theta_0 + \theta_1 x_t + \theta_2 s) + e_{jt}\} \right],$$

where β is a discount factor that we set to 0.9.

Define the ex-ante value functions in period t as the discounted sum of current and future payoffs before the shock e_{jt} is realized and before decision a_t is made, conditional on choosing optimally in every period including t . Denote these ex-ante value functions by $V(x_t, s)$. Further define the conditional value functions $v_j(x, s)$ as the current period payoff of choice j net of e_{jt} :

$$v_j(x, s) = \begin{cases} \beta V(0, s) & j = 0 \\ \theta_0 + \theta_1 x_t + \theta_2 s + \beta V(x + 1, s) & j = 1. \end{cases}$$

Denote by $p_0(x, s)$ the conditional probability of replacing the engine given x and s . Given the distributional assumptions about the shocks, this will be given by

$$p_0(x, s) = \frac{1}{1 + \exp[v_1(x, s) - v_0(x, s)]}.$$

To carry out the simulations, we recursively derive the value functions $v_j(x, s)$ for each possible combination of x , s and t . We then use these to compute the conditional replacement probabilities for the same set of combinations of variables. We generate data for 1000 buses and 2000 time periods. The mileage of each bus is first set to zero in $t = 0$. We then simulate the choices a_t using the conditional replacement probabilities $p_0(x, s)$. Finally, we restrict the generated data to 30 time periods between $t = 1000$ and $t = 1030$. This is to ensure that our data is close to being drawn from a stationary model. Our final dataset consists of types s , mileages x_t and choices a_t for 1000 buses with 30 time period observations each.

7.2. Simulation Results . This section reports the simulation results using the locally robust estimator described in 3.3. To highlight the gain in using the locally robust version of our estimator, we also generate results for the version of our estimator which is suboptimal under continuous state variables (see Section 3.1). We run 1000 simulations with 1000 buses and 30 time periods each. Each round of the simulations proceeds by first generating a dataset as described in Section 7.1. We randomly split this dataset into two samples, \mathcal{N}_1 and \mathcal{N}_2 . We then parameterize $h(a, x)$ and $g(a, x)$ using a third order polynomial in s , x_t and a_t (i.e all individual and interaction terms up to the third order). This implies $k_\phi = k_r = 16$. The choice probabilities η are estimated using a logit model that is a function of the state variables s and x_t , where the same third order polynomial is used as before. Using only observations from the first sample \mathcal{N}_1 and the estimated choice probabilities $\hat{\eta}$, we then estimate the ω parameters using equation 3.3, and the ξ parameters using equation 3.5. Using the observations from the second sample \mathcal{N}_2 , we finally obtain estimates for the θ^* parameters as the solution to the moment equations 3.14. The non-locally robust version of our estimator instead obtains $\hat{\theta}$ using equation 3.9.

Table 1 shows the results. Column (1) reports the true parameters of the model. Columns (2) – (4) report the results for the version of our estimator which is suboptimal under continuous state variables. The results for our locally robust estimator are reported in columns (5) – (7). Column (5) shows that our estimator produces parameter estimates which are closely centered around the true values. The absolute bias after 1000 simulations is less than half of a percent for all three parameters. These results are comparable to those found by Arcidiacono and Miller (2011) in a similar version of the bus engine replacement problem. However, in contrast to their CCP method, our estimator does not exploit a finite dependence property. When comparing the results from our locally robust estimator to the results from the suboptimal estimator in column (2), it can be seen that the bias is up to 60% smaller for the three parameter estimates. However the variance of the locally robust estimator is higher. This is due to the sample splitting employed in the locally robust procedure.

We view the above results as first evidence supporting the workings of our estimators. Further simulations are in progress to test the properties of our estimator in settings with permanent unobserved heterogeneity (see Section 5), or in dynamic discrete games (see Section 6).

8. CONCLUSIONS

We propose a new estimator for DDC models which overcomes previous computational limitations by combining traditional CCP estimation approaches with a TD method from the Reinforcement Learning literature. In making use of simple matrix inversion techniques, our estimator is computationally very cheap and therefore fast. Unlike previous estimation methods, it is able to handle large state spaces in settings where a finite dependence property does not hold. This is of particular importance in settings with continuous state variables where discretization often gives rise to a very high-dimensional state space, or for the estimation of dynamic discrete games. At the same time, our estimator is as efficient as other approaches in simple versions of

TABLE 1. Simulation Results

	DGP (1)	not locally robust			locally robust		
		TDL (2)	bias (3)	bias (%) (4)	TDL (5)	bias (6)	bias (%) (7)
θ_0 (intercept)	2.0	1.9770 (0.1232)	-0.0230	-1.1494	1.9912 (0.1407)	-0.0088	-0.4412
θ_1 (mileage)	-0.15	-0.1492 (0.0050)	0.0008	0.5572	-0.1493 (0.0065)	0.0007	0.4543
θ_2 (bus type)	1.0	1.0056 (0.0847)	0.0056	0.5593	0.9988 (0.1080)	-0.0012	-0.1183

Notes: The table reports results for 1000 simulations. Column (1) shows the true parameter values in the model. Columns (2) and (5) report the mean and standard deviations for the estimated parameters. Columns (2)-(4) are based on the estimation method without correction function, columns (5)-(7) report results for the locally robust estimator. For both methods, biases are reported in absolute terms and as percentage of the parameter values in the data generating process.

the DDC problem. We prove the statistical properties of our estimator and show that it is consistent and converges at parametric rates. Preliminary Monte Carlo simulations using a version of the famous Rust (1987) engine replacement problem confirm these properties in practice.

REFERENCES

- D. Akerberg, X. Chen, J. Hahn, and Z. Liao, “Asymptotic efficiency of semiparametric two-step gmm,” *Review of Economic Studies*, vol. 81, no. 3, pp. 919–943, 2014.
- V. Aguirregabiria and P. Mira, “Swapping the nested fixed point algorithm: A class of estimators for discrete markov decision models,” *Econometrica*, vol. 70, no. 4, pp. 1519–1543, 2002.
- , “Sequential estimation of dynamic discrete games,” *Econometrica*, vol. 75, no. 1, pp. 1–53, 2007.
- , “Dynamic discrete choice structural models: A survey,” *Journal of Econometrics*, vol. 156, no. 1, pp. 38–67, 2010.
- P. Arcidiacono and R. A. Miller, “Conditional choice probability estimation of dynamic discrete choice models with unobserved heterogeneity,” *Econometrica*, vol. 79, no. 6, pp. 1823–1867, 2011.
- P. Arcidiacono, P. Bayer, F. A. Bugni, and J. James, “Approximating high-dimensional dynamic models: Sieve value function iteration,” in *Structural Econometric Models*. Emerald Group Publishing Limited, 2013, pp. 45–95.
- P. Bajari, C. L. Bankard, and J. Levin, “Estimating dynamic models of imperfect competition,” *Econometrica*, vol. 75, no. 5, pp. 1331–1370, 2007.
- A. Benveniste, M. Métivier, and P. Priouret, *Adaptive algorithms and stochastic approximations*. Springer Science & Business Media, 2012, vol. 22.
- V. Chernozhukov, J. C. Escanciano, H. Ichimura, W. K. Newey, and J. M. Robins, “Locally robust semiparametric estimation,” *Working paper*, 2018.

- A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood from incomplete data via the em algorithm,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, 1977.
- G. Hall, G. J. Hitsch, G. Pauletto, and J. Rust, “A comparison of discrete and parametric approximation methods for continuous-state dynamic programming problems,” 2000.
- V. J. Hotz and R. A. Miller, “Conditional choice probabilities and the estimation of dynamic models,” *The Review of Economic Studies*, vol. 60, no. 3, pp. 497–529, 1993.
- V. J. Hotz, R. A. Miller, S. Sanders, and J. Smith, “A simulation estimator for dynamic models of discrete choice,” *The Review of Economic Studies*, vol. 61, no. 2, pp. 265–289, 1994.
- C. Johnson, “Positive definite matrices,” *The American Mathematical Monthly*, vol. 77, no. 3, pp. 259–264, 1970.
- H. Kasahara and K. Shimotsu, “Pseudo-likelihood estimation and bootstrap inference for structural discrete markov decision models,” *Journal of Econometrics*, vol. 146, no. 1, pp. 92–106, 2008.
- , “Sequential estimation of structural models with a fixed point constraint,” *Econometrica*, vol. 80, no. 5, pp. 2303–2319, 2012.
- M. P. Keane and K. I. Wolpin, “The solution and estimation of discrete choice dynamic programming models by simulation and interpolation: Monte carlo evidence,” *the Review of economics and statistics*, pp. 648–672, 1994.
- W. K. Newey, “The asymptotic variance of semiparametric estimators,” *Econometrica*, vol. 62, no. 6, pp. 1349–1382, 1994.
- , “Convergence rates and asymptotic normality for series estimators,” *Journal of econometrics*, vol. 79, no. 1, pp. 147–168, 1997.
- F. Niu, B. Recht, C. Re, and S. J. Wright, “Hogwild!: A lock-free approach to parallelizing stochastic gradient descent (nips 2011),” *Advances in Neural Information Processing Systems* 24, 2011.
- M. Pesendorfer and P. Schmidt-Dengler, “Asymptotic least squares estimators for dynamic games,” *Review of Economic Studies*, vol. 75, no. 3, pp. 901–928, 2008.
- J. Rust, “Optimal replacement of gmc bus engines: An empirical model of harold zurcher,” *Econometrica*, vol. 55, no. 5, pp. 999–1033, 1987.
- V. Semenova, “Machine learning for dynamic models of imperfect information and semiparametric moment inequalities,” *arXiv preprint arXiv:1808.02569*, 2018.
- R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT Press, Cambridge, MA, 2018.
- J. N. Tsitsiklis and B. Van Roy, “An analysis of temporal-difference learning with function approximation,” *IEEE Transactions on Automatic Control*, vol. 42, no. 5, pp. 674–690, 1997.

APPENDIX A. PROOFS OF MAIN RESULTS

In what follows we shall drop the functional argument (a, x) when the context is clear and denote $f' \equiv f(a', x')$ for different functions f .

We start with some useful lemmas:

Lemma 1. *There exists a unique fixed point to the operator $P_\phi \Gamma_z$. If Assumption 1(i) holds, this fixed point is given by $\phi^\top \omega^*$, where ω^* is such that $\mathbb{E}[\phi(z + \beta \phi'^\top \omega^* - \phi^\top \omega^*)] = 0$.*

Proof. First off, we note that Γ_z , and therefore $P_\phi \Gamma_z$, are both contraction maps with the contraction factor β . This implies that $P_\phi \Gamma_z$ has a unique fixed point. Clearly, this fixed point must lie in the space \mathcal{L}_ϕ . Let us denote this as $\phi^\top \omega^*$.

Now for any function $f \in \mathcal{L}_\phi$,

$$\begin{aligned} P_\phi \Gamma_z[f] - f &= \phi^\top \mathbb{E}[\phi \phi^\top]^{-1} \mathbb{E}[\phi(z + \beta f')] - \phi^\top \mathbb{E}[\phi \phi^\top]^{-1} \mathbb{E}[\phi f] \\ &= \phi^\top \mathbb{E}[\phi \phi^\top]^{-1} \mathbb{E}[\phi(z + \beta f' - f)]. \end{aligned}$$

Since $\phi^\top \omega^*$ is the fixed point, we must have

$$\phi^\top \mathbb{E}[\phi \phi^\top]^{-1} \mathbb{E}[\phi(z + \beta \phi'^\top \omega^* - \phi^\top \omega^*)] = 0.$$

But ϕ is linearly independent and $\mathbb{E}[\phi \phi^\top]^{-1}$ is non-singular, by Assumption 1(i). Hence it must be the case

$$\mathbb{E}[\phi(z + \beta \phi'^\top \omega^* - \phi^\top \omega^*)] = 0.$$

This completes the proof the lemma. □

For the next Lemma, we shall use the following definition of a negative-definite matrix: a square, possibly asymmetric, matrix A is said to be negative definite with the coefficient $\bar{\lambda}(A)$ if

$$\sup_{|w|=1} w^\top A w \leq \bar{\lambda}(A) < 0.$$

For a symmetric negative-definite matrix, we have that $\bar{\lambda}(A) = \max \text{eig}(A)$, where $\max \text{eig}(\cdot)$ represents the maximal eigenvalue. We can similarly define a positive definite matrix with the coefficient $\underline{\lambda}(A)$. If the latter is also symmetric, then $\underline{\lambda}(A) = \min \text{eig}(A)$.

We note that under our definition, if A is negative definite, it is also invertible. This holds even if the matrix is asymmetric, see e.g. Johnson (1970).

Lemma 2. *Under Assumption 1(i), the matrix $A := \mathbb{E}[\phi(\beta \phi' - \phi)^\top]$ is negative definite with $\bar{\lambda}(A) \leq -(1 - \beta)\underline{\lambda}(\mathbb{E}[\phi \phi^\top])$, and is therefore invertible.*

Proof. The idea for this proof is taken from Tsitsiklis and van Roy (1997). Recall the definition of $\phi^\top \omega^*$ as the fixed point to $P_\phi T_z[\cdot]$ from Lemma 1. We shall now show that

$$(\omega - \omega^*)^\top A (\omega - \omega^*) \leq -(1 - \beta)\underline{\lambda}(\mathbb{E}[\phi \phi^\top]) |\omega - \omega^*|^2 \quad \forall \omega \in \mathbb{R}^k.$$

Observe that

$$\begin{aligned} A(\omega - \omega^*) &= \mathbb{E} [\phi(z + \beta \phi'^\top \omega - \phi^\top \omega)] - \mathbb{E} [\phi(z + \beta \phi'^\top \omega^* - \phi^\top \omega^*)] \\ &= \mathbb{E} [\phi(z + \beta \phi'^\top \omega - \phi^\top \omega)], \end{aligned}$$

since the second expression in the first equation is 0. Now,

$$\begin{aligned} \mathbb{E} [\phi(z + \beta \phi'^\top \omega - \phi^\top \omega)] &= \mathbb{E} [\phi(a, x)(z(a, x) + \beta \mathbb{E} [\phi(a', x')^\top \omega | a, x] - \phi(a, x)^\top \omega)] \\ &= \mathbb{E} [\phi(\Gamma_z[\phi^\top \omega] - \phi^\top \omega)] \\ &= \mathbb{E} [\phi(P_\phi \Gamma_z[\phi^\top \omega] - \phi^\top \omega)], \end{aligned}$$

where the last equality holds since $\mathbb{E} [\phi(I - P_\phi)\Gamma_z[\phi^\top \omega]] = 0$. We thus have

$$\begin{aligned} (\omega - \omega^*)^\top A(\omega - \omega^*) &= \mathbb{E} [(\omega^\top \phi - \omega^{*\top} \phi)(P_\phi \Gamma_z[\phi^\top \omega] - \phi^\top \omega)] \\ &= \mathbb{E} [(\omega^\top \phi - \omega^{*\top} \phi)(P_\phi \Gamma_z[\phi^\top \omega] - \phi^\top \omega^*)] - \|\phi^\top \omega - \phi^\top \omega^*\|_2^2. \end{aligned}$$

Since $P_\phi \Gamma_z[\cdot]$ is a contraction mapping with contraction factor β , it follows

$$\|P_\phi \Gamma_z[\phi^\top \omega] - \phi^\top \omega^*\|_2^2 = \|P_\phi \Gamma_z[\phi^\top \omega] - P_\phi \Gamma_z[\phi^\top \omega^*]\|_2^2 \leq \beta \|\phi^\top \omega - \phi^\top \omega^*\|_2^2.$$

In view of the above,

$$\begin{aligned} (\omega - \omega^*)^\top A(\omega - \omega^*) &\leq -(1 - \beta) \|\phi^\top \omega - \phi^\top \omega^*\|_2^2 \\ &= -(1 - \beta)(\omega - \omega^*)^\top \mathbb{E}[\phi \phi^\top](\omega - \omega^*) \\ &\leq -(1 - \beta) \underline{\lambda}(\mathbb{E}[\phi \phi^\top]) \|\omega - \omega^*\|^2. \end{aligned}$$

This completes the proof of the lemma. \square

Lemma 3. *Suppose that Assumption 1(i) holds. Then,*

$$\|h - \phi^\top \omega^*\|_2 \leq (1 - \beta)^{-1} \|h - P_\phi[h]\|_2.$$

Proof. Recall that $h(\cdot, \cdot)$ is the unique fixed point of Γ_z , and similarly, $\phi^\top \omega^*$ is the unique fixed point of $P_\phi \Gamma_z$. The operator Γ_z is a contraction mapping with contraction factor β . Furthermore, the projection operator P_ϕ is linear, and $\|P_\phi[f]\|_2 \leq \|f\|_2$ for any function f . Thus

$$\begin{aligned} \|h - \phi^\top \omega^*\|_2 &\leq \|h - P_\phi[h]\|_2 + \|P_\phi[h] - P_\phi \Gamma_z[\phi^\top \omega^*]\|_2 \\ &\leq \|h - P_\phi[h]\|_2 + \|h - \Gamma_z[\phi^\top \omega^*]\|_2 \\ &= \|h - P_\phi[h]\|_2 + \|\Gamma_z[h] - \Gamma_z[\phi^\top \omega^*]\|_2 \\ &\leq \|h - P_\phi[h]\|_2 + \beta \|h - \phi^\top \omega^*\|_2. \end{aligned}$$

Rearranging the above expression proves the desired claim. \square

For the proof of Theorem 1, we shall require some additional notation to take care of the panel dimension when $T > 1$. In this case, while the policy and value functions are stationary, the distribution of (a_{it}, x_{it}) is not time invariant. Let P_t denote the population distribution of (a, x) at time t . Also, let P denote the probability distribution of the process $\{(a_1, x_1), \dots, (a_T, x_T)\}$.

Note that $P \equiv P_1 \times \cdots \times P_T$. We will denote $E[\cdot]$ as the expectation over P . Furthermore, we shall use the $o_p(\cdot)$ and $O_p(\cdot)$ notations to denote convergence in probability, and bounded in probability, respectively, under the probability distribution P .

Note that P is different from \mathbb{P} . The latter provides the distribution of (a, x) after dropping the time index. However, the two are related since for any function f , we can write $\mathbb{E}[f(a, x)] = T^{-1} \sum_{t=1}^T E[f(a_{it}, x_{it})]$ (we could alternatively use this as the definition of $\mathbb{E}[\cdot]$ itself). Note that due to the Markov process assumption, the conditional distribution $P(a_{t+1}, x_{t+1} | a_t, x_t)$ is always independent of t (indeed, one could always consider t as also a part of x). Hence, $\mathbb{P}(a', x' | a, x) \equiv P(a_{t+1}, x_{t+1} | a_t, x_t)$ and $\mathbb{E}[f(a', x') | a, x] \equiv E[f(a_{t+1}, x_{t+1}) | a_t, x_t]$ for all t .

A.1. Proof of Theorem 1. That ω^* exists follows from Lemma 1. To prove that $\hat{\omega}$ exists, it suffices to show that $\hat{A} := \mathbb{E}_n[\phi(\beta\phi' - \phi)^\top]$ is invertible with probability approaching 1. Recall that by our notation above, $\hat{A} = (nT)^{-1} \sum_{i,t} \phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top$, while $A = T^{-1} \sum_t E[\phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top]$. We can thus write $|\hat{A} - A| \leq T^{-1} \sum_t |\hat{A}_t - A_t|$, where $\hat{A}_t := n^{-1} \sum_i \phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top$ and $A_t := E[\phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top]$. Now, by Assumption 1(ii), $|\phi(a, x)|_\infty \leq M$ independent of k_ϕ . We then have

$$\begin{aligned} E|\hat{A}_t - A_t|^2 &= E\left|\frac{1}{n} \sum_i \phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top - E[\phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top]\right|^2 \\ &\leq \frac{1}{nT^2} \sum_i E|\phi_{it}(\beta\phi_{it+1} - \phi_{it})^\top|^2 \leq \frac{k_\phi^2 M^4}{n}. \end{aligned}$$

This proves $|\hat{A}_t - A_t| = O_p(k_\phi/\sqrt{n})$. But T is fixed, which implies that $|\hat{A} - A| = O_p(k_\phi/\sqrt{n})$ as well. We thus obtain $\bar{\lambda}(\hat{A}) \leq \bar{\lambda}(A) + |\hat{A} - A| \leq \bar{\lambda}(A) + o_p(1)$. Since $\bar{\lambda}(A) < 0$, this proves that $\bar{\lambda}(\hat{A}) < 0$ with probability approaching 1, and subsequently, that \hat{A} is invertible. This completes the proof of the first claim.

The second claim follows directly from Lemma 3 and Assumption 1(iii).

For the third claim, let us define $b = \mathbb{E}[\phi z]$ and $\hat{b} = \mathbb{E}_n[\phi z]$. We then have $A\omega^* = b$ and $\hat{A}\hat{\omega} = \hat{b}$. We can combine the two equations to get

$$\hat{A}(\hat{\omega} - \omega^*) = (\hat{b} - b) + (A - \hat{A})\omega^*.$$

The above implies

$$(\hat{\omega} - \omega^*)^\top (-\hat{A})(\hat{\omega} - \omega^*) = (\hat{\omega} - \omega^*)^\top (b - \hat{b}) + (\hat{\omega} - \omega^*)^\top (\hat{A} - A)\omega^*. \quad (\text{A.1})$$

Now, earlier in the proof we have showed that $|\hat{A} - A| = O_p(k_\phi/\sqrt{n})$. Hence it follows $\underline{\lambda}(-\hat{A}) \geq \underline{\lambda}(-A) + o_p(1)$. We thus have

$$(\hat{\omega} - \omega^*)^\top (-\hat{A})(\hat{\omega} - \omega^*) \geq c(1 - \beta)\underline{\lambda}(\mathbb{E}[\phi\phi^\top])|\hat{\omega} - \omega^*|^2, \quad (\text{A.2})$$

with probability approaching 1, for any constant $c \in (0, 1)$. In view of (A.1) and (A.2),

$$|\hat{\omega} - \omega^*| \leq \frac{1}{c(1 - \beta)\underline{\lambda}(\mathbb{E}[\phi\phi^\top])} \left(|\hat{b} - b| + |\hat{A}\omega^* - A\omega^*| \right),$$

with probability approaching 1.

It thus remains to bound $|\hat{b} - b|$ and $|\hat{A}\omega^* - A\omega^*|$. By similar arguments as before, we can define $\hat{b}_t = n^{-1} \sum_i \phi_{it} z_{it}$ and $b_t = E[\phi_{it} z_{it}]$ to obtain

$$E |\hat{b}_t - b_t|^2 = E \left| \frac{1}{n} \sum_i \{\phi_{it} z_{it} - E[\phi_{it} z_{it}]\} \right|^2 \leq \frac{1}{n} E |\phi_{it} z_{it}|^2.$$

This proves

$$E |\hat{b} - b|^2 \leq \sum_t E |\hat{b}_t - b_t|^2 \leq \frac{T}{n} E [|\phi z|^2] \leq \frac{k_\phi T L^2 M^2}{n} = O_p(k_\phi/n).$$

In a similar vein,

$$E |\hat{A}\omega^* - A\omega^*|^2 = E \left| \frac{1}{nT} \sum_{i,t} \{\phi_{it} (\beta \phi_{it+1} - \phi_{it})^\top \omega^* - E[\phi_{it} (\beta \phi_{it+1} - \phi_{it})^\top \omega^*]\} \right|^2 = O_p(k_\phi/n),$$

as long as $E [|\phi (\beta \phi - \phi)^\top \omega^*|^2] = O(k_\phi)$. But the latter is true under Assumptions 1(ii)-(iv) since

$$\mathbb{E} [|\phi (\beta \phi^\top \omega^* - \phi^\top \omega^*)|^2] \leq k_\phi M^2 (2 + 2\beta^2) \mathbb{E} [|\phi^\top \omega^*|^2]$$

and

$$\mathbb{E} [|\phi^\top \omega^*|^2]^{1/2} \leq \|\phi^\top \omega^* - h\|_2 + \|h\|_2 \leq O(k_\phi^{-\alpha}) + (1 - \beta)^{-1} L < \infty,$$

where the second inequality uses the facts $\|\phi^\top \omega^* - h\|_2 = O(k_\phi^{-\alpha})$ (as shown in the second claim of this theorem), and $|h(\cdot, \cdot)|_\infty \leq (1 - \beta)^{-1} |z(\cdot, \cdot)|_\infty < (1 - \beta)^{-1} L$ (which can be easily verified using (2.1) and Assumption 1(iv)). Combining the above, we thus conclude there exists $C < \infty$ such that

$$|\hat{\omega} - \omega^*| \leq C \sqrt{\frac{k_\phi}{n}},$$

with probability approaching one. This completes the proof of the third claim.

Finally, to prove the last claim, observe that

$$\begin{aligned} \|\phi^\top \hat{\omega} - h\|_2^2 &\leq 2 \|\phi^\top \hat{\omega} - \phi^\top \omega^*\|_2^2 + 2 \|\phi^\top \omega^* - h\|_2^2 \\ &= 2(\hat{\omega} - \omega^*)^\top \mathbb{E}[\phi \phi^\top] (\hat{\omega} - \omega^*)^{1/2} + 2 \|\phi^\top \omega^* - h\|_2^2 \\ &\leq \bar{\lambda}(\mathbb{E}[\phi \phi^\top]) O_p\left(\frac{k_\phi}{n}\right) + O_p(k_\phi^{-\alpha}), \end{aligned}$$

where the second inequality follows from the second and third claims of this Theorem. But

$$\bar{\lambda}(\mathbb{E}[\phi \phi^\top]) \leq \|\phi\|_2^2 \leq M^2 k_\phi,$$

by Assumption 1(iv). Combining the above proves the last claim.

A.2. Proof of Theorem 2. We note that the proofs of the first two claims follows from analogous arguments as used in the proof of Theorem 1. We thus only need consider the third claim of the theorem. The fourth claim is a straightforward consequence of this.

Recall that we use a cross-fitting procedure for estimating ξ^* . Let n_1, n_2 denote the sample sizes in the two folds. Also let $\hat{\eta}_1, \hat{\xi}_1$ and $\hat{\eta}_2, \hat{\xi}_2$ denote the estimates of η and ξ^* from the two folds. We shall show that $|\hat{\xi}_1 - \xi| = O_p(\sqrt{k_r/n})$. By a symmetric argument, we will also have

$|\hat{\xi}_2 - \xi| = O_p(\sqrt{k_r/n})$, from which we can conclude $|\hat{\xi} - \xi| = O_p(\sqrt{k_r/n})$. To this end, let $A_r := \mathbb{E}[rr^\top]$, $b_r := \mathbb{E}[r(a, x)e(a', x')]$, $\hat{A}_r^{(1)} := \mathbb{E}_n^{(1)}[rr^\top]$ and $\hat{b}_r^{(1)} := \mathbb{E}_n^{(1)}[r(a, x)e(a', x'; \hat{\eta}_2)]$, where $\mathbb{E}_n^{(1)}[\cdot]$ denotes the empirical expectation using only the observations from the first block. We shall also employ the notation $\psi(a, x, a', x'; \eta) := r(a, x)e(a', x'; \eta)$ and $\psi_{it}(\eta) := r(a_{it}, x_{it})e(a_{it+1}, x_{it+1}; \eta)$.

Based on the above definitions, we have $\hat{A}_r^{(1)}\hat{\xi}_1 = \hat{b}_r^{(1)}$, and $A_r\xi^* = b_r$. Comparing with the proof of Theorem 1, we find that the only difference is in the treatment of $|\hat{b}_r^{(1)} - b_r|$. As in that proof, define $\hat{b}_{rt}^{(1)} := n^{-1} \sum_i \psi_{it}(\hat{\eta}_2)$ and $b_{rt} := E[\psi_{it}(\eta)]$. We then have $|\hat{b}_r^{(1)} - b_r| = T^{-1} \sum_t |\hat{b}_{rt}^{(1)} - b_{rt}|$. Since T is finite, it suffices to bound $|\hat{b}_{rt}^{(1)} - b_{rt}|$ for some arbitrary t . Now, by similar arguments as in the proof of Theorem 1, we have

$$\frac{1}{n_1} \sum_{i=1}^{n_1} \{\psi_{it}(\eta) - E[\psi_{it}(\eta)]\} = O_p\left(\sqrt{k_r/n}\right).$$

Hence the claim follows once we show

$$\hat{b}_{rt}^{(1)} - b_{rt} = \frac{1}{n_1} \sum_{i=1}^{n_1} \{\psi_{it}(\eta) - E[\psi_{it}(\eta)]\} + o_p\left(\sqrt{k_r/n}\right). \quad (\text{A.3})$$

We now prove (A.3). Let \mathcal{N}_2 denote the set of all observations in the second fold. We have

$$\begin{aligned} \hat{b}_{rt}^{(1)} - b_{rt} &= \frac{1}{n_1} \sum_{i=1}^{n_1} \{\psi_{it}(\eta) - E[\psi_{it}(\eta)]\} \\ &= \frac{1}{n_1} \sum_{i=1}^{n_1} \{(\psi_{it}(\hat{\eta}_2) - \psi_{it}(\eta)) - (E[\psi_{it}(\hat{\eta}_2)|\mathcal{N}_2] - E[\psi_{it}(\eta)])\} + \{E[\psi_{it}(\hat{\eta}_2)|\mathcal{N}_2] - E[\psi_{it}(\eta)]\} \\ &:= R_{1nt} + R_{2nt}. \end{aligned}$$

First consider the term R_{1nt} . Define

$$\delta_{it} := (\psi_{it}(\hat{\eta}_2) - \psi_{it}(\eta)) - (E[\psi_{it}(\hat{\eta}_2)|\mathcal{N}_2] - E[\psi_{it}(\eta)]).$$

Clearly, $E[\delta_{it}|\mathcal{N}_2] = 0$. We then have

$$E[|R_{1nt}|^2|\mathcal{N}_2] = \frac{1}{n_1} E[|\delta_{it}|^2|\mathcal{N}_2] = \frac{1}{n_1} E[|\psi_{it}(\hat{\eta}_2) - \psi_{it}(\eta)|^2|\mathcal{N}_2]. \quad (\text{A.4})$$

Now for any (a, x, a', x') , we can note from the definition of $\psi(\cdot)$ that with probability approaching 1,

$$\begin{aligned} |\psi(a, x, a', x'; \hat{\eta}_2) - \psi(a, x, a', x'; \eta)| &\leq |r(a, x)| \{|\ln \hat{\eta}_2 - \ln \eta| + |\hat{\eta}_2 - \eta|\} \\ &\leq M\sqrt{k_r} \{|\ln \hat{\eta}_2 - \ln \eta| + |\hat{\eta}_2 - \eta|\} \\ &\leq M\sqrt{k_r}(2\delta^{-1} + 1)|\hat{\eta}_2 - \eta|, \end{aligned} \quad (\text{A.5})$$

where the second inequality follows from Assumption 2(iii), and the third inequality follows from Assumption 2(v).⁷ Thus in view of (A.4) and (A.5), there exists $C < \infty$ such that

$$E[|R_{1nt}|^2] \leq \frac{Ck_r}{n_1} E[|\hat{\eta}_2(a_{it+1}, x_{it+1}) - \eta(a_{it+1}, x_{it+1})|^2] \leq \frac{Ck_r T}{n_1} \|\hat{\eta}_2 - \eta\|_2^2 = o_p(k_r/n).$$

⁷In particular, we have used the fact $\hat{\eta}_2 > \delta + o_p(1)$ which follows from $\eta > \delta$ and $|\hat{\eta}_2 - \eta| = o_p(1)$.

This proves

$$|R_{1nt}| = o_p(\sqrt{k_r/n}). \quad (\text{A.6})$$

Next consider the term R_{2nt} . We note that $E[\psi_{it}(\eta)]$ is twice Fréchet differentiable. In the main text we have shown that $\partial_\eta E[\psi_{it}(\eta)] = 0$ (c.f equation (3.7)). Furthermore, following some straightforward algebra it is possible to show $|\partial_\eta^2 E[\psi_{it}(\eta)]| \leq C_1 \sqrt{k}$, for some $C_1 < \infty$, as long as η is bounded away from 0 (as assured by Assumption 2(v)). Hence

$$\begin{aligned} E[|R_{2nt}|] &\leq C_1 \sqrt{k_r} E \left[|\hat{\eta}_2(a_{it+1}, x_{it+1}) - \eta(a_{it+1}, x_{it+1})|^2 \right] \\ &\leq C_1 T \sqrt{k_r} \|\hat{\eta}_2 - \eta\|_2^2 = o_p(\sqrt{k_r/n}) \end{aligned} \quad (\text{A.7})$$

where the last equality follows by Assumption 2(v).

Together, (A.6) and (A.7) imply (A.3), which concludes the proof of the theorem.

A.3. Proof of Theorem 4. Let ω_l denote the l th update of ω . From Algorithm 1, we observe that the gradient updates are of the form

$$\omega_{l+1} = \omega_l + \alpha_\omega^{(l)} (z_{it}\phi_{it} - \phi_{it}(\phi_{it} - \beta\phi_{it+1})^\top \omega_l).$$

By standard results on stochastic approximation algorithms (see, e.g, Benveniste et al. (2012), Theorem 17), the above sequence of updates converges to a fixed point $\hat{\omega}$ satisfying

$$\mathbb{E}_n[z\phi - \phi(\phi - \beta\phi')^\top \hat{\omega}] = 0$$

as long as (1) $\mathbb{E}_n[z\phi]$ is finite, (2) $A_n := \mathbb{E}_n[\phi(\phi - \beta\phi')^\top]$ is negative definite, and (3) the learning rate $\alpha_\omega^{(k)}$ satisfies the requirements specified Assumption 3. The first condition is obviously satisfied under Assumption 1(iii). The second condition, that A_n is negative definite, has already been shown in the context of the proof of Theorem 1. Hence, with probability approaching 1, all the three conditions are satisfied and the sequence ω_k converges to $\hat{\omega}$. A similar analysis also applies to gradient descent updates of ξ .