

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

Dynamic Inventory Repositioning in On-Demand Rental Networks

Saif Benjaafar

Department of Industrial and Systems Engineering, University of Minnesota, saif@umn.edu

Daniel Jiang

Department of Industrial Engineering, University of Pittsburgh, drjiang@pitt.edu

Xiang Li

Target Corporation, xiang.li2@target.com

Xiaobo Li

Department of Industrial Systems Engineering and Management, National University of Singapore, iselix@nus.edu.sg

We consider a rental service with a fixed number of rental units distributed across multiple locations. The units are accessed by customers without prior reservation and on an on-demand basis. Customers can decide on how long to keep a unit and where to return it. Because of the randomness in demand and in returns, there is a need to periodically reposition inventory away from some locations and into others. In deciding on how much inventory to reposition and where, the system manager balances potential lost sales with repositioning costs. Although the problem is increasingly common in applications involving on-demand rental services, not much is known about the nature of the optimal policy for systems with a general network structure or about effective approaches to solving the problem. In this paper, first, we show that the optimal policy in each period can be described in terms of a well-specified region over the state space. Within this region, it is optimal not to reposition any inventory, while, outside the region, it is optimal to reposition but only such that the system moves to a new state that is on the boundary of the no-repositioning region. We also provide a simple check for when a state is in the no-repositioning region. Second, we leverage the features of the optimal policy, along with properties of the optimal cost function, to propose a provably convergent approximate dynamic programming algorithm to tackle problems with a large number of dimensions.

Key words: rental networks; inventory repositioning; optimal policies, approximate dynamic programming algorithms, stochastic dual dynamic programming

1. Introduction

We consider a rental service with a fixed number of rental units distributed across multiple locations. Customers are able to rent a unit for one or more periods without a prior reservation and

without specifying, at the time of initiating the rental, the duration of the rental or the return location. That is, customers are allowed to return a unit rented at one location to any other location. We refer to such a service as being *on-demand* and *one-way*. Demand that cannot be fulfilled at the location at which it arises is considered lost and incurs a lost sales penalty (or is fulfilled through other means at an additional cost). Because of the randomness in demand, the length of the rental periods, and return locations, there is a need to periodically reposition inventory away from some locations and into others. Inventory repositioning is costly and the cost depends on both the origin and destination of the repositioning. The service provider is interested in minimizing the lost revenue from unfulfilled demand (lost sales) and the cost incurred from repositioning inventory (repositioning cost). Note that more aggressive inventory repositioning can reduce lost sales but leads to higher repositioning cost. Hence, the firm must carefully mitigate the trade-off between demand fulfillment and inventory repositioning.

Problems with the above features are increasingly common in applications involving on-demand rental services¹. We are particularly motivated by a variety of one-way car sharing services that allow customers to rent from one location and return to another location. Examples include the car sharing service ShareNow (a merger of formerly Car2Go and DriveNow) and similar recently launched services such as LimePod, BlueLA, and BlueSG. These services are on-demand in the sense that they do not require an ahead-of-time reservation or the specification of a return duration ahead of use. They are one-way in the sense that customers are not required to return a vehicle to the same location from which it was rented and instead are allowed to decide on a location, among those in the network, that is most convenient. In other words, these services let customers decide on when and where to return a vehicle, with this information not necessarily shared with the rental service until the rental terminates. For example, LimePod advertises its service as being *free-floating*, offering customers the ability to walk up to any parked vehicle, initiate a rental via its mobile app, and then return the vehicle to any legal parking location within the service region at any time. Customers are charged based on the duration of the rental, with charges calculated at the time the rental terminates.

A challenge in managing these services (and other examples of on-demand rental systems) is the spatial mismatch between vehicle supply and demand that arises from the uncertainty in rental origins, duration, and destinations. Unless adequately mitigated with the periodic repositioning of inventory, the spatial mismatch between supply and demand can lead to a significant loss in revenue.

¹ Renting may become even more prevalent as the economy shifts away from a model built on the exclusive ownership of resources to one based on on-demand access and resource sharing (Sundararajan 2016) (Benjaafar and Hu 2019).

Although the problem is increasingly common in applications involving on-demand rental services², not much is known about the nature of the optimal policy for systems with a general network structure or about effective approaches to solving the problem. This lack of results appears to be due to the multidimensional nature of the problem (i.e., more than one inventory location), compounded by the presence of randomness in demand, rental periods, and return locations, as well as lost sales. In this paper, we address these limitations through two main contributions. The first contribution is theoretical and the second is computational:

- On the theoretical side, we offer a characterization of the optimal policy for the dynamic inventory repositioning problem in a general network setting, accounting for randomness in trip volumes, duration, origin, and destination as well as spatial and temporal dependencies (e.g., likelihood of a trip terminating somewhere being dependent on its origin as well as trip volumes that are dependent on time and location).
- On the computational side, we describe a new cutting-plane-based approximate dynamic programming (ADP) algorithm that can effectively solve the dynamic repositioning problem to near-optimality. We provide a proof of convergence for our algorithm that takes a fundamentally different view from existing cutting-plane-based approaches (this can be viewed as a theoretical contribution to the ADP literature, independent from the repositioning application). We also propose a clustering-based extension to our ADP method that scales to large-scale systems and illustrate its effectiveness on problems with up to 100 locations.

Specifically, we formulate the repositioning problem as a multi-period stochastic dynamic program. We show that the optimal policy in each period can be described in terms of two well-specified regions over the state space. If the system is in a state that falls within one region, it is optimal not to reposition any inventory (we refer to this region as “the no-repositioning” region). If the system is in a state that is outside this region, then it is optimal to reposition some inventory but only such that the system moves to a new state that is on the boundary of the no-repositioning region. Moreover, we provide a simple check for when a state is in the no-repositioning region, which also allows us to compute the optimal policy more efficiently.

One of the distinctive features of the problem considered lies in its non-linear state update function. This non-linearity introduces difficulties in showing the convexity of the problem that must be solved in each period. To address this difficulty, we leverage the fact that the state update function

²Other applications where the periodic repositioning of inventory is important include bike share systems where customers can pick up a bike from one location and return it to any other location within the service region; shipping container rentals in the freight industry where containers can be rented in one location and returned to a different location, with locations corresponding in some cases to ports in different countries; and the use of certain medical equipment, such as IV pumps and wheelchairs, in large hospitals by different departments located in various areas of the hospital.

is piecewise-affine and derive properties for the directional derivatives of the value function. This approach has potential applicability to other systems with piecewise-affine state update functions.

Due to the curse of dimensionality, the optimal policy (and the value function) can be difficult to compute for problems with more than a small number of dimensions. To address this issue, we leverage the results obtained regarding the structure of both the value function and the optimal policy to construct an approximate dynamic programming algorithm. The algorithm combines aspects of approximate value iteration (see, for example, De Farias and Van Roy (2000) and Munos and Szepesvári (2008)) and stochastic dual dynamic programming (see Pereira and Pinto (1991)). We conduct numerical experiments to illustrate the effectiveness of jointly utilizing value function and policy structure, which, to our knowledge, has not yet been explored by related methods in the literature.

The rest of the paper is organized as follows. In Section 2, we review related literature. In Section 3, we describe and formulate the problem. In Sections 4 and 5, we give our structural results for the optimal value function and the optimal policy, respectively. Next, in Section 6 we describe our ADP approach, along with several numerical studies. In Section 7, we address the problem of scaling the ADP algorithm to large systems via the clustering-based extension. In Section 8, we provide concluding comments.

Notation. Throughout the paper, the following notation will be used. We use \mathbf{e} to denote a vector of all ones, \mathbf{e}_i to denote a vector of zeros except 1 at the i th entry, and $\mathbf{0}$ to denote a vector of all zeros (the dimension of these vectors will be clear from the context). Also, we write $\Delta_{n-1}(M)$ to denote the $(n-1)$ -dimensional simplex, i.e., $\Delta_{n-1}(M) = \{(x_1, \dots, x_n) \mid \sum_{i=1}^n x_i = M, \mathbf{x} \geq \mathbf{0}\}$. Similarly, we use $S_n(M)$ to denote the n -dimensional simplex with interior, i.e., $S_n(M) = \{(x_1, \dots, x_n) \mid \sum_{i=1}^n x_i \leq M, \mathbf{x} \geq \mathbf{0}\}$. Throughout, we use ordinary lowercase letters (e.g., x) to denote scalars, and boldfaced lowercase letters (e.g., \mathbf{x}) to denote vectors. The Euclidean norm is denoted $\|\cdot\|_2$. For functions f_1 and f_2 with domain \mathcal{X} , let $\|f_1\|_\infty = \sup_{x \in \mathcal{X}} |f_1(x)|$ and let $f_1 \leq f_2$ denote $f_1(x) \leq f_2(x)$ for all $x \in \mathcal{X}$. We denote the boundary of a set E by $\mathcal{B}(E)$, and the interior of E by E° .

2. Literature Review

There is growing literature on inventory repositioning in car and bike sharing systems; see for example Nair and Miller-Hooks (2011), Shu et al. (2013), Bruglieri et al. (2014), O'Mahony and Shmoys (2015), Freund et al. (2016), Liu et al. (2016), Ghosh et al. (2017), Schuijbroek et al. (2017), Li et al. (2018), Shui and Szeto (2018), Nyotta et al. (2019), and the references therein. Most of this literature focuses on the static repositioning problem, where the goal is to find the optimal placement of vehicles before demand arises, with no more repositioning being made afterwards (e.g.,

repositioning overnight for the next day). The objective function in the associated optimization problem typically accounts for repositioning and user dissatisfaction costs (e.g., lost sales). Much of this work employs mixed integer programming formulations and focuses on the development of algorithms and heuristics. Similarly, the papers that focus on *dynamic* repositioning generally consider heuristic solution techniques and do not offer structural results regarding the optimal policy (see, for example, Ghosh et al. (2017) and Li et al. (2018)).

A notable exception is Li and Tao (2010) who study a finite horizon problem with two locations. They show that the optimal policy in the last period is characterized by a two-limit control policy (a single dimensional version of the policy we show is true in general). They conjecture that the optimal policy has a similar structure in other periods but do not provide a proof. In this paper, we prove that this conjecture is indeed true. Moreover, we prove that a generalized version of the policy holds for problems with more than two locations.

A related stream of literature models vehicle sharing systems as closed queueing networks; see, for example, George and Xia (2011), Wasserhole and Jost (2016), Banerjee et al. (2017), Braverman et al. (2019), Banerjee et al. (2018) and Benjaafar et al. (2021). This literature treats time as being continuous with demand (in the form of an arrival process) that is typically stationary. A common objective in this literature is to identify control policies that maximize a function of the amount of demand satisfied. Control levers include demand throttling (e.g., via pricing), vehicle dispatching (deciding on how to allocate available vehicles to demand as it arises), and empty vehicle repositioning. In general, optimal dynamic policies, such as the ones we consider in this paper, are difficult to characterize. Instead much of this literature relies on analyzing asymptotic regimes, when either the number of vehicles goes to infinity or both the number of vehicles and demand go to infinity. Braverman et al. (2019) consider the optimal repositioning problem in the asymptotic regime where both demand and number of vehicles are allowed to go to infinity. The resulting static repositioning policy is shown to provide an upper bound on the optimal objective function for the finite problem. The static repositioning problem is also discussed in Benjaafar et al. (2021) under a demand balance assumption and using an approximation for vehicle availability at each location.

Wasserhole and Jost (2016) and Banerjee et al. (2017) consider the optimal pricing problem in the asymptotic regime when the number of vehicles goes to infinity. The resulting static pricing policy is shown to provide guaranteed bounds for the finite system. Banerjee et al. (2018) consider the optimal dispatching problem for a stylized setting under the infinite number of vehicles regime and show that the resulting policy provides strong performance bounds.

In this paper, we take a different approach by studying optimal dynamic policies which are generally state-dependent. We consider a setting where time is discretized and allow for demand to be non-stationary. However, we do treat the inventory of vehicles as continuous (this is also the case

under the asymptotic regimes considered in the queueing-based literature). We suspect, though we do not prove it, that the optimal dynamic policy for vehicle repositioning under a queueing model may share similar features to the optimal policy we show for our setting.

Other related papers include Chung et al. (2018) who analyze incentive-based repositioning policies for bike sharing, and Bimpikis et al. (2019) and Balseiro et al. (2020) who consider the spatial aspect of pricing in ride-sharing networks, a related problem to ours, and Ma et al. (2020) who consider a setting with full information to design a spatial-temporal mechanism for prices and wages with desirable properties. Other work considers strategic issues such as fleet sizing, service region design, infrastructure planning, and user dissatisfaction; see, for example, Jian et al. (2016), Raviv and Kolka (2013), He et al. (2017), Lu et al. (2017), Freund et al. (2017), Kabra et al. (2020), and Kaspi et al. (2017). Comprehensive reviews of the literature on vehicle and bike sharing can be found in He et al. (2019) and Freund et al. (2019).

There is literature that addresses inventory repositioning that arises in other settings, including in the context of repositioning of empty containers in the shipping industry, empty railcars in railroad operations, cars in traditional car rentals, and emergency vehicles; see, for example, Lee and Meng (2015) for a comprehensive review. The literature on empty container repositioning is particularly extensive. However that literature focuses on simple networks and relies on heuristics when considering more general problems; see for example Song (2005) and Li et al. (2007). To our knowledge, there are no results regarding the optimal policy for a general network. There is also extensive literature on emergency vehicle repositioning when the demand from different locations is random. Berman (1981) introduces a dynamic programming formulation of the problem. More recent work includes Maxwell et al. (2010) who describe an ADP approach and Maxwell et al. (2014) who provide a lower bound on the performance of repositioning policies; see Bélanger et al. (2019) for a comprehensive review of this literature.

The paper that is closest to ours is He et al. (2020), which was subsequent to an earlier version of this paper³ and which considers a problem similar to ours and solves it using a robust optimization approach. Further discussion of this paper can be found in Section 4. Another subsequent paper that is similar to both our work and He et al. (2020) is Zhao et al. (2020): they derive a structural result for the special case of two locations, two periods, and no ongoing rentals when there is a fixed cost to repositioning. Zhao et al. (2020) build upon the models first introduced in our paper⁴ and the work of He et al. (2020).

³ The first version of our paper appeared online ahead of the first version of He et al. (2020). He et al. (2020) refer to that version of our paper.

⁴ Zhao et al. (2020) cite the working version of this paper.

The problem we consider in this paper shares features with the well-studied dynamic portfolio optimization problem when there are transaction costs. The dynamic portfolio optimization problem involves periodically reallocating funds among different assets, taking into account the stochastic nature of how the value of these assets evolves over time. Constantinides (1979) shows that the structure of the optimal policy resembles that of the optimal policy we describe in this paper for the vehicle sharing problem; see also Leland (1999). That is, it is optimal to do nothing if the system state (defined by the current values of the assets) is within a specified region; otherwise, it is optimal to reallocate funds so that the system state after reallocation lies on the boundary of the do-nothing region. Eberly and Van Mieghem (1997) consider a broader class of resource/capacity allocation problems and prove a similar structure for the optimal policy; see Van Mieghem (2003) for a review of related literature. A discussion of computational approaches, bounds on the optimal solution, and heuristics can be found in Muthuraman and Kumar (2006) and Brown and Smith (2011) and the references therein. The dynamics of the problem we consider are different (e.g., the amount of demand at one location in one period affects the distribution of vehicles at other locations in future periods). Moreover, in our case, the total capacity in the system must be held constant. Our problem is neither a special case nor a generalization of the dynamic portfolio optimization problem. However, these similarities hint that both problems may belong to a more general class of problems whose optimal solution has such a feature.

Finally, there is related literature on computational methods that can solve problems with convex value functions. Some well-known cutting-plane-based approaches are the *stochastic decomposition* algorithm of Higle and Sen (1991), the *stochastic dual dynamic programming* (SDDP) method introduced in Pereira and Pinto (1991), and the *cutting plane and partial sampling* approach of Chen and Powell (1999). Our method is most closely related to SDDP, where full expectations are computed at each iteration. Linowsky and Philpott (2005), Philpott and Guan (2008), Shapiro (2011), and Girardeau et al. (2014) provide convergence analyses of SDDP, but these analyses are designed for finite-horizon problems (or two-stage stochastic programs) and rely on an exact terminal value function and/or that there only exist a finite number of cuts.

Our algorithm is most closely related to the cutting plane methods for the infinite horizon setting proposed in Birge and Zhao (2007) and Warrington et al. (2019). Birge and Zhao (2007) proves uniform convergence of the value function approximations to optimal value for the case of linear dynamics, given a strong condition that the cut in each iteration is computed at a state where a Bellman error criterion is approximately maximized. Computation of such a state is a *difference of convex functions* optimization problem (or a suitable approximation). Warrington et al. (2019) focus on the deterministic setting, use a fixed set of sampled states at which cuts are computed, and do not show consistency of their algorithm. Our algorithm removes these restrictions, yet we

are still able to show uniform convergence to the optimal value function. In particular, our analysis allows for non-linear dynamics and cuts to be computed at states sampled from a distribution. Furthermore, our use of policy structure (i.e., the no-repositioning region characterization) in an SDDP-like algorithm is new.

As an alternative to cutting plane algorithms, Godfrey and Powell (2001) and Powell et al. (2004) propose methods based on stochastic approximation (see Kushner and Yin (2003)) to estimate scalar or separable convex functions, where a piecewise-linear approximation is updated iteratively via noisy samples while ensuring that convexity is maintained. Nascimento and Powell (2009) extend the technique to a finite-horizon ADP setting for the problem of lagged asset acquisition (single inventory state) and provides a convergence analysis; see also Nascimento and Powell (2010). However, these methods are not immediately applicable to our situation, where the value function is multi-dimensional.

3. Problem Formulation

We consider a product rental network consisting of n locations and N rental units. Inventory levels are reviewed periodically and, in each period, a decision is made on how much inventory to reposition away from one location to another. Inventory repositioning is costly and the cost depends on both the origins and destinations of the repositioning. The review periods are of equal length and decisions are made over a specified planning horizon, either finite or infinite.

Demand in each period is positive and random, with each unit of demand requiring the usage of one rental unit for one or more periods, with the rental period being also random. Demand that cannot be satisfied at the location at which it arises is considered lost and incurs a lost sales penalty. A location in the context of a free-floating car sharing system may correspond to a specified geographic area (e.g., a zip code area, a neighborhood, or a set of city blocks). Units rented at one location can be returned to another. Hence, not only are rental durations random but so are return destinations. At any time, a rental unit can be either at one of the locations, available for rent, or in an “ongoing rental” state with a customer.

The sequence of events in each period is as follows. At the beginning of the period, inventory level at each location is observed. A decision is then made on how much inventory to reposition away from one location to another. Subsequently, demand is realized at each location followed by the realization of product returns. Our model assumes, for tractability, that repositioning occurs within a single review period⁵. Note that the solution we obtain is still implementable (feasible) even if the assumption regarding the repositioning time does not hold.

⁵ Similar assumptions on relocation/travel times have been made in much of the existing literature on this topic. He et al. (2020) assume that both customer trips and repositioning trips can be completed within a period. Balseiro et al. (2020) and Waserhole and Jost (2016) assume that travel times are instantaneous. Bimpikis et al. (2019) assumes that going from one location to another takes one period.

We index the periods by $t \in \mathbb{N}$, with $t = 1$ indicating the first period in the planning horizon. We let $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,n}) \in \mathbb{R}^n$ denote the vector of inventory levels before repositioning in period t , where $x_{t,i}$ denotes the corresponding inventory level at location i . Our model uses continuous inventory levels for tractability. This means that the resulting repositioning decisions produced by the model will be continuous. To implement them in a discrete system, one would need to perform a rounding step⁶. Similarly, we let $\mathbf{y}_t = (y_{t,1}, \dots, y_{t,n}) \in \mathbb{R}^n$ denote the vector of inventory levels after repositioning in period t , where $y_{t,i}$ denotes the corresponding inventory level at location i . Note that inventory repositioning should always preserve the total on-hand inventory. Therefore, we require $\sum_{i=1}^n y_{t,i} = \sum_{i=1}^n x_{t,i}$. As we will make clear later, \mathbf{x}_t is only a part of the state in our dynamic system. The second part of the state is the vector of ongoing rentals, defined below.

Inventory repositioning is costly and, for each unit of inventory repositioned away from location i to location j , a cost of c_{ij} is incurred. Consistent with our motivating application of a car sharing system, we assume there is a cost associated with the repositioning of each unit; see He et al. (2020) for similar treatment. Let $\mathbf{c} = (c_{ij})$ denote the cost matrix and let w_{ij} denote the amount of inventory to be repositioned away from location i to location j . Then, the minimum cost associated with repositioning from an inventory level \mathbf{x} to another inventory level \mathbf{y} is given by the solution to the following linear program:

$$\begin{aligned} \min \quad & \mathbf{c} \cdot \mathbf{w} \\ \text{subject to} \quad & \sum_{i=1}^n w_{ij} - \sum_{k=1}^n w_{jk} = y_j - x_j \quad \forall j = 1, \dots, n \\ & \mathbf{w} \geq 0. \end{aligned}$$

The first constraint ensures that the change in inventory level at each location is consistent with the amounts of inventory being moved into ($\sum_i w_{ij}$) and out of ($\sum_k w_{jk}$) that location. The second constraint ensures that the amount of inventory being repositioned away from one location to another is always nonnegative so that the associated cost is accounted for in the objective. It is clear that the value of the linear program depends only on $\mathbf{z} = \mathbf{y} - \mathbf{x}$. Define

$$\begin{aligned} C(\mathbf{z}) = \min \quad & \mathbf{c} \cdot \mathbf{w} \\ \text{subject to} \quad & \sum_{i=1}^n w_{ij} - \sum_{k=1}^n w_{jk} = z_j \quad \forall j = 1, \dots, n \\ & \mathbf{w} \geq 0, \end{aligned} \tag{1}$$

for any $\mathbf{z} \in \mathcal{H}$ where $\mathcal{H} := \{\mathbf{z} \in \mathbb{R}^n : \sum_{i=1}^n z_i = 0\}$. Then the inventory repositioning cost from \mathbf{x} to \mathbf{y} is $C(\mathbf{y} - \mathbf{x})$. Without loss of generality, we assume that $c_{ij} \geq 0$ satisfy the triangle inequality (i.e., $c_{ik} \leq c_{ij} + c_{jk}$ for all i, j, k).

⁶This is reasonable when the number of rental units N is large and is consistent with treatment elsewhere in the literature (see for example He et al. (2020), Li and Tao (2010), and Zhao et al. (2020)) and in much of the literature on stochastic inventory control (see for example Zipkin (2000)).

We let $\mathbf{d}_t = (d_{t,1}, \dots, d_{t,n})$ denote the vector of random demands in period t , with $d_{t,i}$ corresponding to the demand at location i . The amount of demand that cannot be fulfilled is given by $(d_{t,i} - y_{t,i})^+ = \max(0, d_{t,i} - y_{t,i})$. Let β_i denote the per unit lost sales penalty incurred in location i . Then, the total lost sales penalty incurred in period t across all locations is given by $L(\mathbf{y}_t, \mathbf{d}_t) = \sum_{i=1}^n \beta_i (d_{t,i} - y_{t,i})^+$. We assume that each product can be rented at most once within a review period, that is, rental periods are longer than review periods.

To model the randomness in the rental return process, we assume that, at the end of each period t , a random fraction $p_{t,ij}$ of products rented from location i is returned to location j for all $i, j \in \{1, 2, \dots, n\}$, with the rest continuing to be rented. We let \mathbf{P}_t denote the matrix of random fractions, i.e.,

$$\mathbf{P}_t = \begin{pmatrix} p_{t,11} & \cdots & p_{t,1n} \\ \vdots & \ddots & \vdots \\ p_{t,n1} & \cdots & p_{t,nn} \end{pmatrix}.$$

The i th row of \mathbf{P}_t must satisfy $\sum_{j=1}^n p_{t,ij} \leq 1$. The case where $\sum_{j=1}^n p_{t,ij} < 1$ corresponds to a setting where rentals are not immediately returned, while the case where $\sum_{j=1}^n p_{t,ij} = 1$ corresponds to a setting where rental periods are exactly equal to one. Let μ_t denote the joint distribution of \mathbf{d}_t and $\mathbf{P}_t, i = 1, 2, \dots, n$. We assume that the random sequence $(\mathbf{d}_t, \mathbf{P}_t)$ is independent over time, and the expected aggregate demand in each period is finite (i.e., $\int_0^\infty \sum_{i=1}^n d_{t,i} d\mu_t < +\infty$). However, we allow \mathbf{d}_t and \mathbf{P}_t to be dependent. The randomness of \mathbf{P}_t is consistent with the on-demand nature of many rental services, where the provider does not have information regarding the return destination.

Finally, let $\gamma_{t,i}$ for $i = 1, 2, \dots, n$ and $t = 1, 2, \dots, T$ denote the quantity of the product rented from location i that remains outstanding at the beginning of period t . Let $\rho \in [0, 1)$ be the rate at which future costs are discounted.

The model we described above can be formulated as a Markov decision process. Fix a time period t . The system states correspond to the on-hand inventory levels \mathbf{x}_t and the outstanding inventory levels $\boldsymbol{\gamma}_t$. The state space is specified by the $(2n - 1)$ -dimensional simplex, i.e., $(\mathbf{x}_t, \boldsymbol{\gamma}_t) \in \Delta_{2n-1}(N)$. Throughout the paper, we denote $S := S_n(N)$ and $\Delta := \Delta_{2n-1}(N)$ since these notations are frequently used. Actions correspond to the vector of target inventory levels \mathbf{y}_t . Given state $(\mathbf{x}_t, \boldsymbol{\gamma}_t)$, the action space is an $(n - 1)$ -dimensional simplex, i.e., $\mathbf{y}_t \in \Delta_{n-1}(e^T \mathbf{x}_t)$. The transition probabilities are induced by the state update function:

$$\begin{aligned} x_{t+1,i} &= (y_{t,i} - d_{t,i})^+ + \sum_{j=1}^n (\gamma_{t,j} + \min(y_{t,j}, d_{t,j})) p_{t,ji} \quad \forall i = 1, 2, \dots, n, \quad t = 1, 2, \dots, T \\ \gamma_{t+1,i} &= (\gamma_{t,i} + \min(y_{t,i}, d_{t,i})) (1 - \sum_{j=1}^n p_{t,ij}) \quad \forall i = 1, 2, \dots, n, \quad t = 1, 2, \dots, T. \end{aligned}$$

Given a state $(\mathbf{x}_t, \boldsymbol{\gamma}_t)$ and an action \mathbf{y}_t , the repositioning cost is given by $C(\mathbf{y}_t - \mathbf{x}_t)$, and the expected lost sales penalty is given by $l_t(\mathbf{y}_t) = \int L_t(\mathbf{y}_t, \mathbf{d}_t) d\mu_t = \int \sum_i \beta_i (d_{t,i} - y_{t,i})^+ d\mu_t$. The single-period cost is the sum of the inventory repositioning cost and lost sales penalty $r_t(\mathbf{x}_t, \boldsymbol{\gamma}_t, \mathbf{y}_t) =$

$C(\mathbf{y}_t - \mathbf{x}_t) + l_t(\mathbf{y}_t)$. The objective is to minimize the expected discounted cost over a specified planning horizon. In the case of a finite planning horizon with T periods, the optimality equations are given by

$$v_t(\mathbf{x}_t, \gamma_t) = \min_{\mathbf{y}_t \in \Delta_{n-1}(e^T \mathbf{x}_t)} r_t(\mathbf{x}_t, \gamma_t, \mathbf{y}_t) + \rho \int v_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}) d\mu_t \quad (2)$$

for $t = 1, 2, \dots, T$, and $v_{T+1}(\mathbf{x}_{T+1}, \gamma_{T+1}) = 0$ where ρ is the discount factor introduced above.

It is useful to note that the problem to be solved in each period can be expressed in the following form:

$$v_t(\mathbf{x}_t, \gamma_t) = \min_{\mathbf{y}_t \in \Delta_{n-1}(e^T \mathbf{x}_t)} C(\mathbf{y}_t - \mathbf{x}_t) + u_t(\mathbf{y}_t, \gamma_t), \quad (3)$$

where

$$u_t(\mathbf{y}_t, \gamma_t) = \int U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t) d\mu_t, \quad (4)$$

and

$$U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t) = L_t(\mathbf{y}_t, \mathbf{d}_t) + \rho v_{t+1}(\tau_x(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t), \tau_\gamma(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t)), \quad (5)$$

where

$$\tau_x(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) = (\mathbf{y} - \mathbf{d})^+ + \mathbf{P}^T(\gamma + \min(\mathbf{y}, \mathbf{d})), \quad \tau_\gamma(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) = (\gamma + \min(\mathbf{y}, \mathbf{d})) \circ (\mathbf{e} - \mathbf{P}_t \mathbf{e}), \quad (6)$$

where \circ denotes the Hadamard product (or the entrywise product), i.e., $(a_1, a_2, \dots, a_n) \circ (b_1, b_2, \dots, b_n) = (a_1 b_1, a_2 b_2, \dots, a_n b_n)$. The next two assumptions state some useful conditions on the return fractions \mathbf{P}_t and the repositioning costs c_{ij} .

ASSUMPTION 1. *Let $p_{\min} \in (0, 1]$ be a constant. For every period t , there exists a random variable $p_t \in [p_{\min}, 1]$ such that $\sum_{j=1}^n p_{t,ij} = \sum_{j=1}^n p_{t,kj} = p_t$, $\forall i, k = 1, 2, \dots, n$. An equivalent statement is that $p_{t,ij} = p_t \tilde{q}_{t,ij}$ for some $\tilde{q}_{t,ij}$ where $\sum_{j=1}^n \tilde{q}_{t,ij} = 1$ for all i .*

Assumption 1 implies that the probability of a vehicle being returned in a given period does not depend on the location at which the vehicle is rented, but the distribution of the return locations does depend on the origin⁷

ASSUMPTION 2. *The repositioning costs satisfy $\rho c_{\max} - c_{\min} \leq p_{\min} (\beta_i - c_{\min})$ for all $i = 1, \dots, n$, where $c_{\max} = \max_{i,j} c_{ij}$ and $c_{\min} = \min_{i,j; i \neq j} c_{ij}$.*

⁷ We make this assumption for tractability of the theoretical analysis in Section 4, but note that it is plausible because whether to return and where to return are usually two separate decisions for customers. Furthermore, in the case of rental networks located in dense urban regions, we expect many rental locations to have similar properties in terms of customers' rental/return behaviors. In Appendix A.1 we provide empirical support for this assumption based on real data obtained from the one-way car sharing service Car2Go.

The second assumption enforces boundedness in the difference of cost parameters, with the upper bound depending on p_{\min} . If $p_{\min} = 1$, where the rental duration is always one period (corresponding to the setting of He et al. (2020)), the restriction reduces to $\rho c_{\max} \leq \beta_i$ for all i . This means that the cost of lost sales outweighs the cost of inventory repositioning in the next period⁸. If $p_{\min} < 1$, the assumption prevents the unpleasant situation where one might want to deliberately “hide” the inventory due to the difference in the repositioning cost. It is clear from the assumption that $\rho c_{\max} - c_{\min} \leq p_{\min}(\beta_i - c_{\min}) \leq p_t(\beta_i - c_{\min})$ for all i .

Under Assumptions 1 and 2, we are able to show (see Section 4 and 5) that the value function in each period, consisting of the lost sales and the cost-to-go as defined next, is always convex, which is perhaps surprising given the non-linear state update and the lost sales feature.

4. Convexity of $u_t(\mathbf{y}_t, \boldsymbol{\gamma}_t)$

The main purpose of this section is to establish the convexity of $u_t(\mathbf{y}_t, \boldsymbol{\gamma}_t)$ defined in (5) for all periods t . We will also show that a similar result holds for the infinite-horizon case. These results will allow us later on (see Section 5) to characterize the structure of the optimal policy. They will also be useful in developing an efficient solution procedure (see Section 6).

4.1. The Finite-Horizon Problem

In this section, we consider the finite-horizon problem discussed in section 3. For the last period (period T), $u_T(\mathbf{y}_T, \boldsymbol{\gamma}_T) = l_T(\mathbf{y}_T)$ is clearly convex. A natural question is whether the convexity of $u_t(\cdot)$ is preserved when we consider previous periods. The main difficulty is that the state update in (6) is non-linear. However, if we introduce an auxiliary variable $\boldsymbol{\omega} = \min\{\mathbf{d}, \mathbf{y}\}$, the state update function could be written in the following linear form:

$$\tau_x(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}, \mathbf{P}) = \mathbf{y} - \boldsymbol{\omega} + \mathbf{P}^T(\boldsymbol{\gamma} + \boldsymbol{\omega}), \quad \tau_\gamma(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}, \mathbf{P}) = (\boldsymbol{\gamma} + \boldsymbol{\omega}) \circ (\mathbf{e} - \mathbf{P}_t \mathbf{e}).$$

We show that we can replace the constraint $\boldsymbol{\omega} = \min\{\mathbf{d}, \mathbf{y}\}$ with $\boldsymbol{\omega} \leq \min\{\mathbf{d}, \mathbf{y}\}$, which would then imply the convexity of $u_t(\cdot)$ (see Appendix A.3 for the analysis, given as a series of technical lemmas). Let

$$u'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}, \boldsymbol{\eta}) = \lim_{t \downarrow 0} \frac{u(\mathbf{x} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) - u(\mathbf{x}, \boldsymbol{\gamma})}{t} \quad (7)$$

denote the directional derivative of $u(\cdot)$ at $(\mathbf{x}, \boldsymbol{\gamma})$ along the direction $(\mathbf{z}, \boldsymbol{\eta})$. We call $(\mathbf{z}, \boldsymbol{\eta})$ a *feasible direction* at $(\mathbf{x}, \boldsymbol{\gamma})$ if $(\mathbf{x} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) \in \Delta$ for small enough $t > 0$. The main results are summarized in the following theorem.

⁸ A similar but slightly weaker condition is assumed in He et al. (2020). In this sense, the convexity in our paper can include their results as a special case except that in their model, lost sales costs depend on both origin and destination and the return destinations are known by the platform at the time of rental. In our case, we assume, consistent with the reality of many one-way vehicle sharing systems, that the destination of a rental is not revealed until a realized trip is completed. In settings where the lost sales cost depends on both the origin and destination of a trip, β_i has the interpretation of the expected lost sales cost over all destinations.

THEOREM 1. *Suppose Assumptions 1 and 2 hold. For $t = 1, \dots, T$, both $u_t(\cdot)$ defined in (4) and $v_t(\cdot)$ defined in (2) are convex and continuous in Δ . Moreover, the following properties of $u_t(\mathbf{y}_t, \gamma_t)$ hold for $t = 1, \dots, T$:*

(a) $u_t(\mathbf{y}_t, \gamma_t) = \mathbb{E}_{\mathbf{d}_t, \mathbf{P}_t} [U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t)]$ where $U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t)$ can be reformulated as the following convex optimization program

$$\begin{aligned}
 U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t) = \min_{\boldsymbol{\omega}, \mathbf{x}_{t+1}, \gamma_{t+1}} & \sum_{i=1}^n \beta_i (d_{t,i} - \omega_i) + \rho v_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}) \\
 \text{subject to} & \mathbf{x}_{t+1} = \mathbf{y}_t - \boldsymbol{\omega} + \mathbf{P}_t^T (\gamma_t + \boldsymbol{\omega}), \\
 & \gamma_{t+1} = (\gamma_t + \boldsymbol{\omega}) \circ (\mathbf{e} - \mathbf{P}_t \mathbf{e}), \\
 & \boldsymbol{\omega} \leq \mathbf{y}_t, \text{ and} \\
 & \boldsymbol{\omega} \leq \mathbf{d}_t;
 \end{aligned} \tag{8}$$

(b) $|u'_t(\mathbf{y}_t, \gamma_t; \pm \boldsymbol{\eta}, \mp \boldsymbol{\eta})| \leq \sum_{i=1}^n \beta_i \eta_i$ for all $(\mathbf{y}_t, \gamma_t) \in \Delta$ and any feasible direction $(\pm \boldsymbol{\eta}, \mp \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$;

(c) $u'_t(\mathbf{y}_t, \gamma_t; \mathbf{0}, \mathbf{z}) \leq (\rho/2) c_{\max} \sum_{i=1}^n |z_i|$ for all $(\mathbf{x}, \gamma) \in \Delta$ and any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = 0$;

(d) $u_t(\cdot)$ is Lipschitz continuous on Δ° with Lipschitz constant $(3/2) \sqrt{2n} \beta_{\max}$, where $\beta_{\max} = \max_i \beta_i$.

A comprehensive proof of Theorem 1 can be found in Appendix A.3. Here, we give an outline of the approach. We apply induction, starting from $v_{T+1}(\mathbf{y}, \gamma) = 0$. We show in Proposition 3 and Proposition 4 that if $v_{t+1}(\cdot)$ is convex and satisfies certain bounds on its directional derivatives, then for any realization of $\mathbf{d}_t, \mathbf{P}_t$, the function $U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t)$ can be reformulated as the convex program (8) and satisfies two types of bounds on its directional derivatives. The first type (item 1) shows that if we turn some of the available inventory into ongoing rentals, the reduced or enlarged cost can be upper bounded by the lost sales cost of these products. The same bound holds if we remove some of the ongoing rentals and make them available at the locations from which they were rented. The second type bound states that if we change the origin of some of the ongoing rentals (i.e., we change γ only), the difference in cost can be upper bounded by the product of $(\rho c_{\max}/2)$ and the one-norm of the difference in γ . The primary reason is that the total return fraction for period t , p_t , does not depend on the origin. Therefore, the difference of costs is at most the repositioning cost in the next period. To complete the induction, we show in Proposition 5 that given the convexity of $u_t(\mathbf{y}_t, \gamma_t)$ and aforementioned bounds on its directional derivatives, $v_t(\mathbf{x}_t, \gamma_t)$ is convex and satisfies the directional derivative bounds required by Proposition 3 and 4. Finally, we show that item (b) and (c) imply item (d).

We note that although the formulation in (8) has similarities to the result in Lemma 1 in He et al. (2020), our proof technique is fundamentally different. First, to show that the reformulation is exact, we need to show that if $\boldsymbol{\omega} \neq \min(\mathbf{d}_t, \mathbf{y}_t)$, we can increase some components of $\boldsymbol{\omega}$ so that the

objective function is not worse while keeping it feasible. For the problem in He et al. (2020) (which corresponds to our problem with $p_{\min} = 1$), since all the outstanding cars are returned at the end of the period, we can arrive at any state \mathbf{x}_{t+1} by adjusting the decision variable y_t . However, when $p_{\min} < 1$, a change in ω would result in a change in γ_{t+1} , which can not be rebalanced by changing y_t . To show the monotonicity of ω , we need the directional derivative of $v_{t+1}(\cdot)$ to satisfy delicate bounds as required in Proposition 3. Second, though the aforementioned bound of $v_{t+1}(\cdot)$ is clearly true for the last period, for general t , we require that the directional derivatives of $u_t(\cdot)$ satisfy the two types of bounds in Theorem 1. To show that these bounds indeed hold, we carry out careful convex analysis and induction per Proposition 3, Proposition 4 and Proposition 5. Above all, it is highly non-trivial to show the exactness of the reformulation in (8). Our proof technique might be of independent interest for high-dimensional inventory problems with lost sales.

An immediate consequence of the reformulation is that, if all random variables satisfy discrete distributions, the problem can be written as a large-scale linear program.

COROLLARY 1. *Suppose Assumptions 1 and 2 hold and $(\mathbf{d}_t, \mathbf{P}_t)$ follows a discrete distribution for all t , then the optimal policy $\pi^* = (\pi_1^*, \dots, \pi_T^*)$ can be computed as the optimal solution to a large-scale linear program.*

Corollary 1 allows us to approximate the optimal policy by replacing each expectation with the finite sum of a few samples and solve the large-scale LP. However, since the size of the LP grows exponentially in the number of samples and the number of periods, for reasonable values of T , we can only afford to solve the problem with a single sample path. In Section 6, we use one sample (the mean demand) for each period to approximate the optimal policy and compare it with our ADP solution procedure.

4.2. The Infinite-Horizon Problem

We have shown that $u_t(\cdot)$ is convex for each period for the finite-horizon problem. Next we show that the same can be said about the stationary problem with infinitely many periods. In such a problem, we denote the common distribution for $(\mathbf{d}_t, \mathbf{P}_t)$ by μ . Similarly, we denote the common values of $L_t(\cdot)$, $l_t(\cdot)$ and $r_t(\cdot)$ by $L(\cdot)$, $l(\cdot)$ and $r(\cdot)$, respectively. We use π to denote a stationary policy that uses the same decision rule π in each period. Under π , the state of the process is a Markov random sequence $\{(X_t, \Gamma_t), t = 1, 2, \dots\}$. The optimization problem can be written as a Markov decision process (MDP):

$$v(\mathbf{x}, \gamma) = \min_{\pi} \mathbb{E}_{\mathbf{x}}^{\pi} \left\{ \sum_{t=1}^{\infty} \rho^{t-1} r(X_t, \Gamma_t, \pi(X_t, \Gamma_t)) \right\}, \quad (9)$$

where $X_1 = \mathbf{x}$ a.e. is the initial state of the process. Let $\tilde{v}_T(\mathbf{x}, \gamma) = \min_{\pi} \mathbb{E}_{\mathbf{x}}^{\pi} \left\{ \sum_{t=1}^T \rho^{t-1} r(X_t, \Gamma_t, \pi_t(X_t, \Gamma_t)) \right\}$ denote the value function of a stationary problem with

T periods. It is well known that the functions $\tilde{v}_T(\cdot)$ converge uniformly to $v(\cdot)$ and $v(\cdot)$ is the unique solution to

$$v(\mathbf{x}, \gamma) = \min_{\mathbf{y} \in \Delta_{n-1}(\mathbf{e}^T \mathbf{x})} r(\mathbf{x}, \gamma, \mathbf{y}) + \rho \int v(\tau_x(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}), \tau_\gamma(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P})) d\mu, \quad (10)$$

where $\tau_x(\cdot)$ and $\tau_\gamma(\cdot)$ correspond to the state update functions defined in (6), i.e.,

$$\begin{aligned} \tau_x(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) &= (\mathbf{y} - \mathbf{d})^+ + \mathbf{P}^T(\gamma + \min(\mathbf{y}, \mathbf{d})) \quad \forall t = 1, 2, \dots, T, \text{ and} \\ \tau_\gamma(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) &= (\gamma + \min(\mathbf{y}, \mathbf{d})) \circ (\mathbf{e} - \mathbf{P}_t \mathbf{e}) \quad \forall t = 1, 2, \dots, T. \end{aligned} \quad (11)$$

For details, the reader may refer to Chapter 6 of Puterman (1994).

As in the finite-horizon version, the problem to be solved can be written in the following form:

$$v(\mathbf{x}, \gamma) = \min_{\mathbf{y} \in \Delta_{n-1}(\mathbf{e}^T \mathbf{x})} C(\mathbf{y} - \mathbf{x}) + u(\mathbf{y}, \gamma), \quad (12)$$

where

$$u(\mathbf{y}, \gamma) = \int U(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) d\mu, \quad (13)$$

and

$$U(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) = L(\mathbf{y}, \mathbf{d}) + \rho v(\tau_x(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}), \tau_\gamma(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P})). \quad (14)$$

THEOREM 2. *Suppose Assumptions 1 and 2 hold. Both $u(\cdot)$ defined in (13) and $v(\cdot)$ defined in (12) are convex and continuous in Δ . Moreover, we have:*

(a) $U(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P})$ defined in (14) can be reformulated as the following convex optimization program

$$\begin{aligned} U(\mathbf{y}, \gamma, \mathbf{d}, \mathbf{P}) &= \min_{\omega, \tau_x, \tau_\gamma} \sum_{i=1}^n \beta_i (d_i - \omega_i) + \rho v(\tau_x, \tau_\gamma) \\ \text{subject to } \tau_x &= \mathbf{y} - \omega + \mathbf{P}^T(\gamma + \omega), \\ \tau_\gamma &= (\gamma + \omega) \circ (\mathbf{e} - \mathbf{P}\mathbf{e}), \\ \omega &\leq \mathbf{y}, \text{ and} \\ \omega &\leq \mathbf{d}; \end{aligned} \quad (15)$$

(b) $|u'(\mathbf{y}, \gamma; \mp \boldsymbol{\eta}, \pm \boldsymbol{\eta})| \leq \sum_{i=1}^n \beta_i \eta_i$ for all $(\mathbf{x}, \gamma) \in \Delta$ and any feasible direction $(\mp \boldsymbol{\eta}, \pm \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$;

(c) $u'(\mathbf{y}, \gamma; \mathbf{0}, \mathbf{z}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |z_i|$ for all $(\mathbf{x}, \gamma) \in \Delta$ and any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = 0$; and

(d) $u(\cdot)$ is Lipschitz continuous on Δ° with Lipschitz constant $(3/2)\sqrt{2n}\beta_{\max}$, where $\beta_{\max} = \max_i \beta_i$.

5. The Optimal Repositioning Policy

In this section, we characterize the structure of the optimal policy. We do so for both the finite and infinite horizon cases. Recall that, for both cases the repositioning problem can be stated as

$$v(\mathbf{x}, \gamma) = \min_{\mathbf{y} \in \Delta_{n-1}(\mathbf{e}^T \mathbf{x})} C(\mathbf{y} - \mathbf{x}) + u(\mathbf{y}, \gamma) \text{ for } (\mathbf{x}, \gamma) \in \Delta, \quad (16)$$

where $C(\cdot)$ is the repositioning cost specified by (1) and $u(\cdot)$ is a convex and continuous function that maps Δ to $\mathbb{R} \cup \{-\infty, \infty\}$. The principle result of this section is the characterization of the optimal policy through the *no-repositioning set*, the collection of inventory levels from which no repositioning should be made. The no-repositioning set for a function $u(\cdot)$ when the outstanding inventory level is γ can be defined as follows:

$$\Omega_u(\gamma) = \{\mathbf{x} \in \Delta_{n-1}(I) : u(\mathbf{x}, \gamma) \leq C(\mathbf{y} - \mathbf{x}) + u(\mathbf{y}, \gamma) \forall \mathbf{y} \in \Delta_{n-1}(I)\}, \forall \gamma \in S \quad (17)$$

where $I = N - \sum_{i=1}^n \gamma_i$. Note that I is a function of γ (or equivalently \mathbf{x}). For notational simplicity, we suppress the dependency of I on γ (or \mathbf{x}). By definition, no repositioning should be made from inventory levels inside $\Omega_u(\gamma)$. In the following theorem, we show that $\Omega_u(\gamma)$ is non-empty, connected and compact and, for inventory levels outside $\Omega_u(\gamma)$, it is optimal to reposition to some point on the boundary of $\Omega_u(\gamma)$. Recall that we denote the boundary of a set E by $\mathcal{B}(E)$, and the interior of E by E° .

THEOREM 3. *The no-repositioning set $\Omega_u(\gamma)$ is nonempty, connected and compact for all $\gamma \in S$. An optimal policy π^* to (16) satisfies*

$$\begin{aligned} \pi^*(\mathbf{x}, \gamma) &= \mathbf{x} && \text{if } \mathbf{x} \in \Omega_u(\gamma); \\ \pi^*(\mathbf{x}, \gamma) &\in \mathcal{B}(\Omega_u(\gamma)) && \text{otherwise.} \end{aligned} \quad (18)$$

Solving a nondifferentiable convex program such as (16) usually involves some computational effort. One way to reduce this effort, suggested by Theorem 3, is to characterize the no-repositioning set $\Omega_u(\gamma)$. Characterizing the no-repositioning region can help us identify when a state is inside $\Omega_u(\gamma)$, which allows our ADP algorithm to more easily compute the value iteration step; see Section 6. Recall that $u'(\mathbf{x}, \gamma; \mathbf{z}, \boldsymbol{\eta}) = \lim_{t \downarrow 0} \frac{u(\mathbf{x} + t\mathbf{z}, \gamma + t\boldsymbol{\eta}) - u(\mathbf{x}, \gamma)}{t}$ denotes the directional derivative of $u(\cdot)$ at (\mathbf{x}, γ) along the direction $(\mathbf{z}, \boldsymbol{\eta})$. Since $u(\cdot)$ is assumed to be convex and continuous in Δ , $u'(\mathbf{x}, \gamma; \mathbf{z}, \boldsymbol{\eta})$ is well defined for $(\mathbf{x}, \gamma) \in \Delta$. Recall also that $(\mathbf{z}, \boldsymbol{\eta})$ is a *feasible direction* at (\mathbf{x}, γ) if $(\mathbf{x} + t\mathbf{z}, \gamma + t\boldsymbol{\eta}) \in \Delta$ for small enough $t > 0$. In what follows, we provide a series of first order characterizations of $\Omega_u(\gamma)$, the first of which relies on the directional derivatives.

PROPOSITION 1. *$\mathbf{x} \in \Omega_u(\gamma)$ if and only if*

$$u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0}) \geq -C(\mathbf{z}) \quad (19)$$

for any feasible direction $(\mathbf{z}, \mathbf{0})$ at (\mathbf{x}, γ) .

Proposition 1 is essential for several subsequent results. However, using Proposition 1 to verify whether a point lies inside the no-repositioning set is computationally impractical, as it involves checking an infinite number of inequalities in the form of (19). In the following proposition, we provide a second characterization of $\Omega_u(\gamma)$ using the subdifferentials. Before we proceed, we introduce

the following notations. \mathbf{g} is said to be a subgradient of $u(\cdot, \gamma)$ at \mathbf{x} if $u(\mathbf{y}, \gamma) \geq u(\mathbf{x}, \gamma) + \mathbf{g}^T(\mathbf{y} - \mathbf{x})$ for all \mathbf{y} . The set of all subgradients of $u(\cdot, \gamma)$ at \mathbf{x} is denoted by $\partial_x u(\mathbf{x}, \gamma)$. It is well known that $\partial_x u(\mathbf{x}, \gamma)$ is nonempty, closed and convex for \mathbf{x} in the interior, which is equivalent to $\mathbf{x} > \mathbf{0}$ in our setting.

PROPOSITION 2. $\mathbf{x} \in \Omega_u(\gamma)$ if $\partial_x u(\mathbf{x}, \gamma) \cap \mathcal{G} \neq \emptyset$, where $\mathcal{G} = \{(g_1, \dots, g_n) : g_i - g_j \leq c_{ij} \forall i, j\}$. If $\mathbf{x} > \mathbf{0}$, then the converse is also true.

Proposition 2 suggests whether a point lies inside the no-repositioning set depends on whether $u(\cdot, \gamma)$ has certain subgradients at this point. Such a characterization is useful if we can compute the subdifferential $\partial_x u(\mathbf{x}, \gamma)$. In particular, if $u(\cdot, \gamma)$ is differentiable at \mathbf{x} , then $\partial_x u(\mathbf{x}, \gamma)$ consists of a single point $\nabla_x u(\mathbf{x}, \gamma)$. In this case, determining its optimality only involves checking $n(n-1)$ inequalities.

COROLLARY 2. Suppose $u(\cdot, \gamma)$ is differentiable at $\mathbf{x} \in \Delta_{n-1}(I)$. Then, $\mathbf{x} \in \Omega_u(\gamma)$ if and only if $\frac{\partial u(\mathbf{x}, \gamma)}{\partial x_i} - \frac{\partial u(\mathbf{x}, \gamma)}{\partial x_j} \leq c_{ij}$ for all i, j .

The no-repositioning set $\Omega_u(\gamma)$ can take on many forms. We first discuss the case where there are only two locations. In this case, the no-repositioning set corresponds to a closed line segment with the boundary being the two end points. The optimal policy reduces to a state-dependent two-threshold policy.

COROLLARY 3. Suppose $n = 2$. For $\gamma \in S$, let $I = N - \gamma_1 - \gamma_2$. Then $\Omega_u(\gamma) = \{(x, I - x) : x \in [s_1(\gamma), s_2(\gamma)]\}$, where $s_1(\gamma) = \inf\{x : u'((x, I - x, \gamma_1, \gamma_2); (1, -1, 0, 0)) \geq -c_{21}\}$ and $s_2(\gamma) = \sup\{x : -u'((x, I - x, \gamma_1, \gamma_2); (-1, 1, 0, 0)) \leq c_{12}\}$. An optimal policy π^* to (16) satisfies

$$\begin{aligned} \pi^*(x, I - x, \gamma_1, \gamma_2) &= (s_1(\gamma), I - s_1(\gamma)) \text{ if } x < s_1(\gamma), \\ \pi^*(x, I - x, \gamma_1, \gamma_2) &= (x, I - x) \text{ if } s_1(\gamma) \leq x < s_2(\gamma), \\ \pi^*(x, I - x, \gamma_1, \gamma_2) &= (s_2(\gamma), I - s_2(\gamma)) \text{ otherwise.} \end{aligned}$$

Corollary 3 is a direct consequence of Theorem 3, Proposition 1, and the fact that there are only two feasible directions. It shows that the optimal policy to problem (16) in the two-dimensional case is described by two thresholds $s_1(\gamma) < s_2(\gamma)$ on the on-hand inventory level x at location 1. If x is lower than s_1 , it is optimal to bring the inventory level up to s_1 by repositioning inventory from location 2 to location 1. On the other hand, if x is greater than s_2 , it is optimal to bring the inventory level at location 1 down to s_2 . When x falls between s_1 and s_2 , it is optimal not to reposition as the benefit of inventory repositioning cannot offset the cost.

When there are more than two locations, a threshold policy is not naturally defined due to the total inventory constraint. In what follows, we characterize the no-repositioning set for two important special cases, the first of which corresponds to when $u(\cdot, \gamma)$ is a convex quadratic function.

If the demands are uniformly distributed, then for the last period, $u(\cdot, \gamma)$ is a quadratic function since only the lost sales cost is involved. In this case, the no-repositioning set is a polyhedron defined by $n(n-1)$ linear inequalities.

EXAMPLE 1. For a fixed γ , suppose $u(\mathbf{y}, \gamma) = \mathbf{y}^T \mathbf{B}(\gamma) \mathbf{y} + \mathbf{y}^T \mathbf{b}(\gamma) + b_0(\gamma)$ and $\mathbf{B}(\gamma)$ is positive semidefinite. By Corollary 2, $\Omega_u(\gamma) = \{\mathbf{y} \in \Delta_{n-1}(I) : 2\mathbf{y}^T B_i(\gamma) + b_i(\gamma) - 2\mathbf{y}^T B_j(\gamma) - b_j(\gamma) \leq c_{ij} \forall i, j\}$, where $B_i(\gamma)$ is the i -th row of $\mathbf{B}(\gamma)$.

We point out that, in general, the no-repositioning set can be non-convex. The following example illustrates that even if $u(\cdot)$ is smooth, $\Omega_u(\gamma)$ might still be non-convex.

EXAMPLE 2. Suppose $\gamma = \mathbf{0}$, $u(\mathbf{y}) = y_1^3 + y_2^2 + y_3^2$, and $c_{ij} = 0.5$ (note that the inventory state \mathbf{y} is always nonnegative so u is convex). Then, the no-repositioning set is characterized by $\Omega_u = \{\mathbf{y} \in \Delta_{n-1} : -0.5 \leq 3y_1^2 - 2y_3 \leq 0.5, -0.5 \leq 3y_1^2 - 2y_2 \leq 0.5, -0.5 \leq 2y_2 - 2y_3 \leq 0.5\}$.

Note that in Example 2, $u(\mathbf{y})$ is a convex function but the no-repositioning set is not convex because the region under the parabolas $2y_2 - 3y_1^2 = 0.5$ and $2y_3 - 3y_1^2 = 0.5$ is not convex. See Figure 1 for the case where $N = y_1 + y_2 + y_3 = 1$.

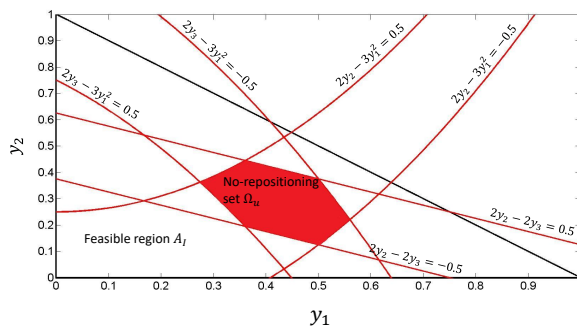


Figure 1 An illustration of a non-convex no-repositioning set

6. An Approximate Dynamic Programming Approach

So far, we have studied the theoretical properties of the repositioning problem. In this section, we propose an approximate dynamic programming algorithm, to which we refer as “R-ADP,” that exploits the structure of both the value function and the optimal policy under a sampled demand and return model. Although Theorems 1 and 2 allow for the use of convex optimization to help resolve the issue of a multi-dimensional continuous action space, the difficulty of a multi-dimensional and continuous state space remains. We refer readers to Bertsekas and Tsitsiklis (1996) and Powell (2007) for a detailed discussion of the computational challenges and solution methods associated with large MDPs. In particular, we note that when the problem size (the number of locations or the number of time periods) is large, simple approximations of continuous problems, such as discretization or aggregation, will usually fail. In addition, discretization can cause our

structural properties to break down, which means the convexity result and characterization of the optimal policy given Theorems 1 and 2 can no longer be readily used. Informal numerical experiments show that even if we do not consider the ongoing rentals (rental period is always one), approximating the dynamic program via discretization within a reasonable accuracy is already a formidable task for a three-location problem.

It is thus necessary for us to consider more scalable techniques. A key feature of the algorithm we describe next is that each iteration involves solving one or more linear programs, allowing it to leverage the scalability and computational advantages of off-the-shelf solvers. We show via numerical experiments that the algorithm can produce high quality solutions on problems with states up to 19 dimensions (10 locations) within a reasonable amount of time. The algorithm also possesses the important theoretical property of asymptotically optimal value function approximations; see Theorem 4. In the rest of this section, we motivate and describe the algorithm, prove its convergence, discuss some practical considerations, and present the numerical results.

6.1. The R-ADP Algorithm

Theorems 1 and 2 describe the most important feature of our dynamic program, that $u(\cdot)$, the summation of current period lost sales and the cost-to-go, is convex and continuous. Moreover, Proposition 2 provides a characterization of when it is optimal not to reposition. Our algorithm takes advantage of these two structural results. It is well known that a convex function can be written as the point-wise supremum of its tangent hyperplanes, i.e.,

$$u(\mathbf{y}, \gamma) = \sup_{\hat{\mathbf{y}}, \hat{\gamma}} u(\hat{\mathbf{y}}, \hat{\gamma}) + (\mathbf{y} - \hat{\mathbf{y}})^T \nabla_{\mathbf{y}} u(\hat{\mathbf{y}}, \hat{\gamma}) + (\gamma - \hat{\gamma})^T \nabla_{\gamma} u(\hat{\mathbf{y}}, \hat{\gamma}).$$

This suggests that we can build an approximation to $u(\cdot)$ by iteratively adding lower-bounding hyperplanes, with the hope that the approximation becomes arbitrarily good when enough hyperplanes are considered. This is the main idea of the algorithm, with special considerations made to account for the complicated structure of the state update functions. Using a lower, piecewise-affine approximation of a convex function is a commonly-used idea in stochastic programming; see, for example, Figure 1 of Philpott and Guan (2008) for an illustration.

Our algorithm is motivated by various aspects of approximate value iteration (see De Farias and Van Roy (2000) and Munos and Szepesvári (2008)) and stochastic dual dynamic programming (see Pereira and Pinto (1991)). The features and analysis that distinguish our ADP algorithm from previous work in the literature are summarized below.

1. Our algorithm has the ability to skip the optimization step when a sampled state is detected as being in the no-repositioning region. This step uses Proposition 2 and it is applied to the value function approximation at every iteration.

2. The underlying model of SDDP and other cutting-plane methods (see, e.g., Hige and Sen (1991), Pereira and Pinto (1991), Birge and Zhao (2007)) is typically a two-stage or multi-stage stochastic *linear* program. In our case, we have non-linear state updates, which makes the optimization step of the algorithm difficult. To sidestep this difficulty, in our algorithm, we approximate $u(\cdot)$ instead of $v(\cdot)$, while computing the state updates outside of the optimization step.

3. Our algorithm is designed for the infinite horizon setting, where each approximation “bootstraps” from the previous approximation and convergence is achieved despite the absence of a terminal condition such as “ $v_{T+1} \equiv 0$ ” used in the finite-horizon case. As such, the convergence analyses used in Chen and Powell (1999), Linowsky and Philpott (2005), Philpott and Guan (2008), Shapiro (2011), and Girardeau et al. (2014) do not apply.⁹ Moreover, we remove a strong condition used in a previous convergence result by Birge and Zhao (2007) for the infinite horizon setting, where cuts are computed at states that approximately maximize a Bellman error criterion. Selecting such a state requires solving a difference of convex functions optimization problem. Our algorithm and proof technique do not require this costly step.

Throughout this section, suppose that we are given M samples of the demand and the return fraction matrix $(\mathbf{d}_1, \mathbf{P}_1), (\mathbf{d}_2, \mathbf{P}_2), \dots, (\mathbf{d}_M, \mathbf{P}_M)$. Our goal is to optimize the sampled model. The idea is to start with an initial piecewise-affine lower approximation $u_0(\mathbf{y}, \boldsymbol{\gamma})$ (such as $u_0(\mathbf{y}, \boldsymbol{\gamma}) = 0$) and then dynamically add linear functions (referred to as cuts in our discussion) into consideration. Suppose we currently have $u_J(\mathbf{y}, \boldsymbol{\gamma}) = \max_{k=1, \dots, N_J} g_k(\mathbf{y}, \boldsymbol{\gamma})$ where $g_k(\mathbf{y}, \boldsymbol{\gamma}) = (\mathbf{y} - \mathbf{y}_k)^T \mathbf{a}_k + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_k)^T \mathbf{b}_k + \iota_k$, and N_J is the total number of cuts in the approximation after iteration J . We then need to evaluate the functional value and the gradient of the following function: $\bar{u}_J(\mathbf{y}, \boldsymbol{\gamma}) = \frac{1}{M} \sum_{s=1}^M \{L(\mathbf{y}, \mathbf{d}_s) + \rho \bar{v}_J(\tau_{\mathbf{x}}(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}_s, \mathbf{P}_s), \tau_{\boldsymbol{\gamma}}(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}_s, \mathbf{P}_s))\}$, where $\bar{v}_J(\mathbf{x}, \boldsymbol{\zeta}) = \min_{\mathbf{z} \in \Delta_{n-1}(\mathbf{e}^T \mathbf{x})} C(\mathbf{z} - \mathbf{x}) + u_J(\mathbf{z}, \boldsymbol{\zeta})$ at a sample point $(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\gamma}})$. Note that $\bar{v}_J(\mathbf{x}, \boldsymbol{\zeta})$ is a linear program. To find out the derivatives of $\bar{v}_J(\mathbf{x}, \boldsymbol{\zeta})$, we write down the dual formulation for $\bar{v}_J(\mathbf{x}, \boldsymbol{\zeta})$ as follows:

$$\begin{aligned} \bar{v}_J(\mathbf{x}, \boldsymbol{\zeta}) = \max & (\lambda_0 \mathbf{e} + \boldsymbol{\lambda})^T \mathbf{x} + \sum_{k=1}^J \mu_k (-\mathbf{a}_k^T \mathbf{y}_k + \mathbf{b}_k^T (\boldsymbol{\zeta} - \boldsymbol{\gamma}_k) + \iota_k) \\ \text{s.t.} & \sum_{k=1}^J \mu_k = 1, \\ & \lambda_i - \lambda_j \leq c_{ij}, \quad \forall i, j = 1, 2, \dots, n, \\ & -\lambda_i + \sum_{k=1}^J \mu_k a_{ki} - \lambda_0 \geq 0, \quad \forall i = 1, 2, \dots, n, \\ & \mu_k \geq 0, \quad \forall k = 1, 2, \dots, N_J. \end{aligned} \quad (20)$$

From (20), we understand that $\nabla_{\mathbf{x}} \bar{v}_J(\mathbf{x}, \boldsymbol{\zeta}) = \lambda_0^* \mathbf{e} + \boldsymbol{\lambda}^*$ and $\nabla_{\boldsymbol{\zeta}} \bar{v}_J(\mathbf{x}, \boldsymbol{\zeta}) = \sum_{k=1}^J \mu_k^* \mathbf{b}_k$, where $(\lambda_0^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ is an optimal solution for problem (20). The Jacobian matrix for the state update function is

$$\nabla_{\bar{\mathbf{x}}, \mathbf{y}} = \text{Diag}(1_{\mathbf{y}_t > \mathbf{d}_k}) + \mathbf{P}_k \circ (1_{\mathbf{y}_t \leq \mathbf{d}_k} \mathbf{e}^T), \nabla_{\bar{\mathbf{x}}, \boldsymbol{\gamma}} = \mathbf{P}_k,$$

⁹ For example, we do not make use of a property that there are only a finite number of distinct cuts; see Lemma 1 of Philpott and Guan (2008). We remark, however, that our algorithm has a natural adaptation for finite-horizon problems.

Parameters: $l, c_{ij}, \rho > 0$
 Data: $\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_M, \mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_M$
 Input: Initial approximation $u_0(\mathbf{y}, \boldsymbol{\gamma}) = \max_{k=1, \dots, N_0} g_k(\mathbf{y}, \boldsymbol{\gamma})$
for $J = 0, 1, 2, \dots$
 Sample a finite set of states \mathcal{S}_{J+1} from Δ according to some distribution
 for $s = 1, 2, \dots, |\mathcal{S}_{J+1}|$
 Let $(\tilde{\mathbf{y}}_s, \tilde{\boldsymbol{\gamma}}_s)$ be the s th sampled state in \mathcal{S}_{J+1}
 Run Subroutine 6 in Appendix A.5 with input $(\tilde{\mathbf{y}}_s, \tilde{\boldsymbol{\gamma}}_s)$ and let the result be $g_{s+N_J}(\mathbf{y}, \boldsymbol{\gamma})$
 end for
 Set $u_{J+1}(\mathbf{y}, \boldsymbol{\gamma}) = \max_{k=1, \dots, N_{J+1}} g_k(\mathbf{y}, \boldsymbol{\gamma})$, where $N_{J+1} = N_J + |\mathcal{S}_{J+1}|$
end for

Table 1 R-ADP Algorithm

$$\nabla_{\bar{\boldsymbol{\gamma}}, \boldsymbol{\gamma}} = \text{Diag}(\mathbf{e} - \mathbf{P}_k \mathbf{e}), \text{ and } \nabla_{\bar{\boldsymbol{\gamma}}, \mathbf{y}} = \text{Diag}((\mathbf{e} - \mathbf{P}_k \mathbf{e}) \circ \mathbf{1}_{\mathbf{y}_t \leq \mathbf{d}_k}),$$

where \bar{x} and $\bar{\boldsymbol{\gamma}}$ stand for $\tau_{\mathbf{x}}$ and $\tau_{\boldsymbol{\gamma}}$ respectively. By computing equation (20) for all pairs of $(\tau_{\mathbf{x}}(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}_i, \mathbf{P}_i), \tau_{\boldsymbol{\gamma}}(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}_i, \mathbf{P}_i))$, we can find the tangent hyperplane of $\bar{u}_J(\mathbf{y}, \boldsymbol{\gamma})$ at $(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\gamma}})$.

While (20) can always be solved, we can apply Proposition 2, a characterization of the no-repositioning region, to reduce the computational load. We first define some terms for $u_J(\mathbf{y}, \boldsymbol{\gamma})$. Let $\mathcal{K} = \{k \mid a_{ki} - a_{kj} \leq c_{ij} \quad \forall i, j\}$ denote the set of cuts that satisfy the no-reposition condition. We also let

$$\begin{aligned} D^k = \{(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta \mid & (\mathbf{y} - \mathbf{y}_k)^T \mathbf{a}_k + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_k)^T \mathbf{b}_k + \boldsymbol{\nu}_k \\ & \geq (\mathbf{y} - \mathbf{y}_l)^T \mathbf{a}_l + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_l)^T \mathbf{b}_l + c_l \quad \forall l = 1, 2, \dots, K\} \end{aligned} \quad (21)$$

denote a subset of the feasible region that is dominated by the k -th cut. Then we have the following lemma.

LEMMA 1. *If $\mathbf{x} \in D^k$ with $k \in \mathcal{K}$, we have $\mathbf{x} \in \Omega_{u_J}(\boldsymbol{\zeta})^{10}$ and one optimal solution for problem (20) is $\boldsymbol{\lambda} = \mathbf{a}_k, \mu_k = 1, \mu_l = 0, \forall l \neq k$.*

The R-ADP algorithm is described in Table 1 while the procedure for adding new cuts is described in Table 6 in Appendix A.5. The essential idea is to iterate the following steps: (1) sample a set of states, (2) compute the appropriate supporting hyperplanes at each state, and (3) add the hyperplanes to the convex approximation of $u(\mathbf{y}, \boldsymbol{\gamma})$. If $u_J(\mathbf{y}, \boldsymbol{\gamma}) \leq u(\mathbf{y}, \boldsymbol{\gamma})$, we have $\bar{u}_J(\mathbf{y}, \boldsymbol{\gamma}) \leq u(\mathbf{y}, \boldsymbol{\gamma})$. Therefore, $g_{s+N_J}(\mathbf{y}, \boldsymbol{\gamma})$, a tangent hyperplane for $\bar{u}_J(\mathbf{y}, \boldsymbol{\gamma})$, is a lower bound for $u(\mathbf{y}, \boldsymbol{\gamma})$, which means that $u_{J+1}(\mathbf{y}, \boldsymbol{\gamma})$ is also a lower bound for $u(\mathbf{y}, \boldsymbol{\gamma})$. Through the course of R-ADP, we obtain an improving sequence of lower approximations to the true $u(\mathbf{y}, \boldsymbol{\gamma})$ function. Hence, if u_0 is a uniform underestimate of u , we know that $u_J(\mathbf{y}, \boldsymbol{\gamma})$ is a bounded monotone sequence and thus its limit exists.

¹⁰ Note that, in general, $\Omega_u(\boldsymbol{\gamma}) \neq \bigcup_{k \in \mathcal{K}} D^k$. The reason is that even if two cuts are both not in \mathcal{K} , the intersection of these two cuts could still include the subgradient that satisfies the no-reposition condition.

There are several reasonable strategies for sampling the set \mathcal{S}_{J+1} . The easiest way is to set $|\mathcal{S}_J| = 1$ (i.e., only add a single cut¹¹ per iteration) and then sample one state according to some distribution over Δ — this is the approach taken in the numerical experiments of this paper. Our implementation of R-ADP also uses an iteration-dependent state sampling distribution to improve the practical performance (see Appendix A.6); therefore, we introduce the following assumption to support the convergence analysis.

ASSUMPTION 3. *On any iteration J , the sampling distribution produces a set \mathcal{S}_J of states from Δ° . The sampled sets $\{\mathcal{S}_J\}_{J=1}^\infty$ satisfy $\sum_{J=1}^\infty \mathbf{P}(\mathcal{S}_J \cap A \neq \emptyset) = \infty$ for any set $A \subseteq \Delta^\circ$ with positive volume.*

This is not a particularly restrictive assumption and should be interpreted simply as requiring an adequate exploration of the state space, a common requirement for ADP and reinforcement learning algorithms (Bertsekas and Tsitsiklis 1996). As an example, for the case of one sample per iteration, one might consider the following sampling strategy, parameterized by a deterministic sequence $\{\epsilon_J\}$: with probability $1 - \epsilon_J$, choose the state in any manner and with probability ϵ_J , select a state uniformly at random over Δ° . In this case, we have that $\mathbf{P}(\mathcal{S}_J \cap A \neq \emptyset) \geq \epsilon_J \cdot \text{volume}(A)$. As long as $\sum_J \epsilon_J = \infty$, Assumption 3 is satisfied.

6.2. Convergence of the R-ADP Algorithm

We are now ready to discuss the convergence of the R-ADP algorithm. For simplicity, we consider the case where $|\mathcal{S}_{J+1}| = 1$ for all iterations J . The extension to the batch case, $|\mathcal{S}_{J+1}| > 1$, follows the same idea and is merely a matter of more complicated notation (note, however, that we will nevertheless make use of a simple special case of batch algorithm as an analysis tool within the proof).

THEOREM 4. *Suppose Assumptions 1, 2, and 3 hold and that R-ADP samples one state per iteration. Suppose the initial value function approximation u_0 is a lower bound on the optimal value function u , and satisfies properties (b) and (c) stated in Theorem 1, namely that*

- $|u'_0(\mathbf{y}, \boldsymbol{\gamma}; \pm \boldsymbol{\eta}, \mp \boldsymbol{\eta})| \leq \sum_{i=1}^n \beta_i \eta_i$ for all $(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta$ and any feasible direction $(\pm \boldsymbol{\eta}, \mp \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$; and
- $u'_0(\mathbf{y}, \boldsymbol{\gamma}; \mathbf{0}, \mathbf{z}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |z_i|$ for all $(\mathbf{x}, \boldsymbol{\gamma}) \in \Delta$ and any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = 0$.

Then, the sequence $\{u_J\}$ converges uniformly and almost surely to the optimal value function u , i.e., it holds that $\|u_J - u\|_\infty \rightarrow 0$ almost surely.

¹¹ If parallel computing is available, one might consider the “batch” version of the algorithm (i.e., $|\mathcal{S}_{J+1}| > 1$) by performing the inner for-loop of Algorithm 1 on multiple processors (or workers). In this case, each worker receives u_J , samples a state, and computes the appropriate supporting hyperplane. The main processor would then aggregate the results into u_{J+1} and start the next iteration by broadcasting u_{J+1} to each worker.

The proof of Theorem 4 relies on relating each sample path of the algorithm to an auxiliary algorithm where the cuts are added in “batches” rather than one by one. We show that, after accounting for the different timescales, the value function approximations generated by R-ADP are close to the approximations generated by the auxiliary algorithm. By noticing that the auxiliary algorithm is an approximate value iteration algorithm whose per-iteration error can be bounded in $\|\cdot\|_\infty$ due to Lemma 8, we quantify its error against exact value iteration, which in turn allows us to quantify the error between R-ADP and exact value iteration. We make use of ϵ -covers of the state space (for arbitrarily small ϵ) along with Assumption 3 to argue that this error converges to zero. Note that one can satisfy the conditions for u_0 by taking it to be a constant function that is a lower bound of u ; for example, $u_0(\cdot) = 0$ is a suitable choice. In Appendix A.6, we discuss two practical aspects associated with implementing R-ADP: (1) checking for and removing redundant cuts and (2) specifying an effective state-sampling distribution.

6.3. Benchmarking R-ADP on Random Problem Instances

We first present some benchmarking results of running R-ADP on a set of randomly generated problems ranging from $n = 2$ to $n = 10$ locations, the largest of which corresponds to a dynamic program with a 19-dimensional continuous state space. We set the discount factor as $\rho = 0.95$, the repositioning costs to be $c_{\min} = c_{\max} = 1$, and the lost sales cost as $\beta_i = 2$. We consider normalized total inventory of $N = 1$, and for each problem instance, we take $M = 50$ demand and return probability samples as follows. With each location i , we associate a truncated normal demand distribution (so that it is nonnegative) with mean ν_i and standard deviation σ_i . The ν_i are drawn from a uniform distribution and then normalized so that $\sum_i \nu_i = 0.3$. We then set $\sigma_i = \nu_i$ so that locations with higher mean demand are also more volatile. Next, we follow Assumption 1 and sample one outcome of a matrix (\tilde{q}_{ij}) such that each row is chosen uniformly from a standard simplex. Each of the M samples of the return probability matrix consists of (\tilde{q}_{ij}) multiplied by a random scaling factor drawn from Uniform(0.7,0.9). Hence, we have $p_{\min} = 0.7$. We compare the performance of R-ADP policy to the performance of several baselines approaches:

- *Myopic Policy.* The myopic policy (Myo.) minimizes the single-period lost sales and repositioning costs, i.e., the policy associated with $v(\cdot) = 0$.
- *Rolling-Horizon Deterministic Lookahead Policy (Mean).* This policy considers a k -period rolling horizon lookahead, obtained by solving the large-scale LP described in Corollary 1 taking $(\mathbf{d}_t, \mathbf{P}_t)$ to be their means. We use the abbreviation ‘ k -RH-M’ to refer to this policy.
- *No-Repositioning Policy.* The no-repositioning policy (No-R) does not reposition any inventory.

We use a maximum of 1000 cuts for all problem instances and we run the R-ADP algorithm for 10,000 iterations for $n \leq 6$ and for 20,000 iterations for $n = 7, 8, 9, 10$. We initially sample

80% of states randomly¹² and 20% of states from the replay buffer of the myopic policy. As the algorithm progresses, we transition toward a distribution of 20% randomly, 0% from the myopic replay buffer, and 80% from the current ADP replay buffer. Note that Assumption 3 is satisfied for this sampling scheme. Redundancy checks are performed every 250 iterations. The performance of the ADP algorithm is evaluated using Monte-Carlo simulation over 500 sample paths (across 20 initial states, randomly sampled subject to zero outstanding rentals) at various times during the training process. Since the ADP algorithm itself is random, we repeat the training process 10 times for each problem instance in order to obtain confidence intervals (which are shown in Figure 2).

n	Sec./Iter.	R-ADP	Myo.	3-RH-M	5-RH-M	7-RH-M	10-RH-M
2	0.06	99.2%	60.1%	64.9%	65.8%	66.0%	65.7%
3	0.21	98.7%	70.7%	71.8%	73.9%	73.5%	75.2%
4	0.27	95.9%	78.1%	69.2%	69.1%	70.1%	69.1%
5	0.22	96.4%	72.5%	70.7%	71.6%	72.9%	73.3%
6	0.29	94.1%	75.3%	74.4%	75.5%	76.1%	76.7%
7	0.42	88.1%	74.2%	59.3%	60.3%	61.3%	61.9%
8	0.44	85.0%	66.0%	61.6%	62.4%	62.1%	63.3%
9	0.48	88.2%	62.2%	57.1%	57.5%	58.5%	58.0%
10	0.54	83.4%	60.8%	49.7%	51.0%	51.6%	52.7%
avg.	-	92.1%	68.9%	64.3%	65.2%	65.8%	66.2%

Table 2 Performance Comparison of Repositioning Policies

The results¹³ are summarized in Table 2. The first column ‘ n ’ shows the number of locations (note that $2n - 1$ is the dimension of the state space). The second column ‘Sec./Iter.’ shows the CPU time for training the R-ADP policy on a 4 GHz Intel Core i7 processor using four cores, which includes the time needed to remove cuts and generate the replay buffer. Figure 3 shows the amount of computational savings when using the no-repositioning region structure (Lemma 1 and Proposition 2), with 95% confidence intervals. We roughly attain 5%-8% CPU savings by making use of the policy structure derived. The remaining columns give a percentage optimality relative to the ADP lower bound for each of the policies, computed as the percentage of the lower bound (LB) achieved when the baseline no-repositioning policy is set as “0% optimal.” (Note that this is a *lower bound* on the percentage optimality relative to the optimal policy). This is done via:

$$\% \text{ optimality to lower bound} = \frac{\text{cost of No-R policy} - \text{cost of R-ADP policy}}{\text{cost of No-R policy} - \text{highest lower bound}}. \quad (22)$$

In terms of wall clock time, we observe that our ADP algorithm produces near-optimal results for $n \leq 6$ within an hour (for $n = 6$, we are using $0.29 \cdot 10000$ seconds or 48 minutes). For the larger

¹² Each sampled state is given by $(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\gamma}}) = (\xi \mathbf{y}', (1 - \xi) \boldsymbol{\gamma}') \in \Delta$ where \mathbf{y}' and $\boldsymbol{\gamma}'$ are independent uniform samples from $\Delta_{n-1}(N)$ and $\xi \sim \text{Uniform}(p_{\min}(M), p_{\max}(M))$ where $p_{\min}(M)$ and $p_{\max}(M)$ are the minimum and maximum row sums of the return fraction matrix over the M samples. This sampling scheme can be considered a nearly uniform sample over the state space, except with the two parts of the state re-scaled by relevant problem parameters so that they are more likely to fall in important regions.

¹³ The same random seed is used in all instances (i.e., all n) to generate the problem parameters.

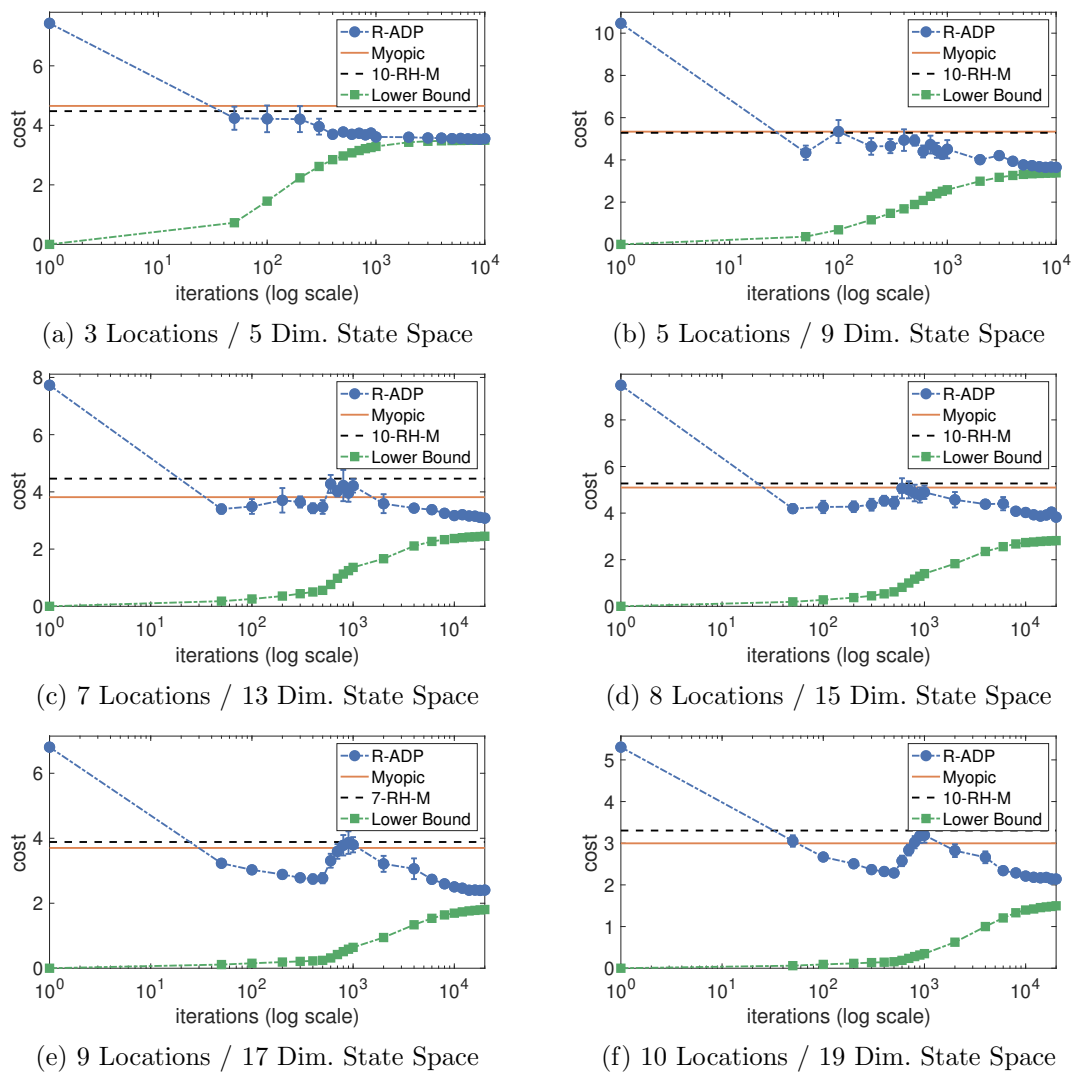


Figure 2 Performance of R-ADP on Randomly Generated Problems

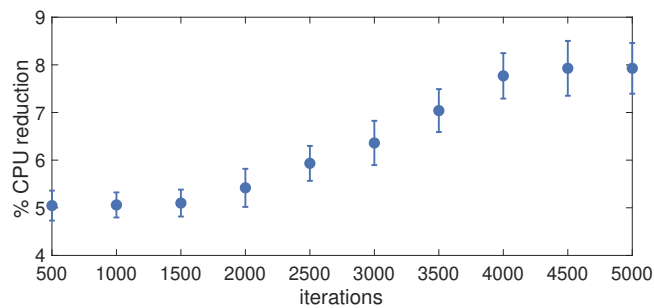


Figure 3 Computational Savings using Policy Structure ($n = 5$)

problems of $n \geq 7$, when provided a limited amount of computation — around three hours for 20,000 iterations — the estimated optimality gap is slightly larger, between 12%–17%. Note that

these are *offline* computation times; the learned policy is implemented by solving a linear program (in roughly the same amount of time as a single iteration of R-ADP, typically less than one second).

Figure 2 shows the performance of the R-ADP policy as the algorithm progresses, along with 95% confidence intervals and lower bounds. We also show the best performing variant of the rolling-horizon deterministic lookahead approach and the myopic policy. In all cases, the cost of the R-ADP policy eventually becomes lower than that of the baselines and approaches the accompanying lower bound. Note that the minor upticks in cost around iteration 1,000 appear to be due to the value function approximation hitting the limit of 1,000 cuts for the first time.

Based on the average performance across the different n (last row of Table 1), the myopic policy comes closest to matching the performance of R-ADP. We note that the total amount of inventory repositioned by the R-ADP policy is considerably higher than that of the myopic policy, between 37% higher (for $n = 5$) and 79% (for $n = 9$) higher. This suggests that the improvement upon the myopic policy can be attributed to a more aggressive repositioning strategy. Since the myopic policy does not take into account customers' return behaviors, the additional repositioning activity observed in the ADP policy can be explained by its attempt to plan for the future by counteracting the effects of \mathbf{P} . In Appendix A.7, we vary parameter settings and provide comparative statics regarding the impact of (1) total demand, (2) demand volatility, (3) return fraction (i.e., fraction of vehicles returned per period), and (4) uniformity of return locations.

7. Scaling-up to Large Systems via CR-ADP

In Section 6.3, we provided computational results for R-ADP on 19-dimensional MDPs ($n = 10$ locations). Approximating MDPs of larger dimensions is well-known to be challenging due to the curse of dimensionality. For example, Lu et al. (2017) approximate a 9-location problem using a two-stage stochastic integer program and He et al. (2020) approximate a 5-location problem using a robust approach within an MDP model without a convergence guarantee. However, practical instances of the repositioning problem may involve much larger values of n . In this section, we show that a surprisingly simple extension of R-ADP via a clustering approach allows it to scale to very large systems (for example, systems with $n = 100$ locations). The approach outperforms the rolling-horizon deterministic lookahead baseline, a commonly used “scalable” strategy for large-scale MDPs (Powell (2007)). Additionally, we show that this method also produces good results when the effective horizon¹⁴ of the MDP is long (i.e., when $\rho = 0.99$).

¹⁴ One way to determine an appropriate discount factor is to connect it with the *effective planning horizon* of the MDP, commonly taken to be $\mathcal{O}(1/(1-\rho))$; see, for example, Jiang et al. (2015).

7.1. CR-ADP

We propose the following simple framework for applying R-ADP to problems with a large number of locations n . We execute the R-ADP algorithm on an auxiliary MDP obtained by clustering locations together so that there is a manageable number of them, and then heuristically deconstruct the “clustered policy” into a policy for the original MDP. This policy is referred to as CR-ADP. Specifically, the algorithm consists of the following steps:

1. First, we partition the n locations into \tilde{n} clusters. Let $\mathcal{C}_k \subseteq \{1, 2, \dots, n\}$ be the k -th cluster for $k = 1, 2, \dots, \tilde{n}$.
2. Using these clusters, we define a transformation of the problem parameters from the n -location instance to the \tilde{n} -location instance. Denote the transformed demand, return fraction, repositioning costs, and lost sales cost by $\tilde{\mathbf{d}}_t$, $\tilde{\mathbf{P}}_t$, $\tilde{\mathbf{c}}$, and $\tilde{\beta}_i$, respectively. We then solve the more tractable \tilde{n} -location problem using R-ADP and obtain a policy $\tilde{\pi}$.
3. To implement the policy: given an n -location state, we obtain an \tilde{n} -location state by summing inventory in each cluster \mathcal{C}_k and produce an \tilde{n} -location repositioning-decision $\tilde{\mathbf{y}}$ using $\tilde{\pi}$. We then use a “splitting” heuristic to construct a n -location repositioning-decision \mathbf{y} from $\tilde{\mathbf{y}}$.

There are many reasonable ways to design the clusters, transform the problem parameters, and to split inventory. We show results for a straightforward instantiation of this framework, where adjacent locations are clustered, demand is summed up within clusters, repositioning costs and return fractions are appropriately averaged, and inventory is split according to the demand proportion of each location relative to the cluster’s total demand. Details are given in Appendix A.8.

Tables 3 and 4 show the results of the clustered approach CR-ADP when compared to the same baseline policies used in Section 6.3, for $\rho = 0.95$ and $\rho = 0.99$. Note that these results are given in terms of expected cost, rather than an optimality percentage, since we do not obtain lower bounds of the original MDP when using CR-ADP. We consider problems with locations ranging from $n = 20$ to $n = 100$, with $\tilde{n} = 10$ clusters in each case. Note that the baselines do not use clustering; in fact, for large values of n , the LP used in the rolling-horizon deterministic lookahead approach becomes computationally intractable for larger values of the lookahead horizon. We see that despite solving a clustered, approximate problem, the ADP approach outperforms all of the baseline policies.

These results point to the benefit of jointly considering downstream values and the stochasticity of demand/return fractions, as neither of these features alone is able to produce high-performing policies (as evidenced by the myopic and deterministic lookahead approaches). Lastly, we remark that many aspects of our clustering heuristic design could be further refined and the most appropriate strategy might be highly problem-dependent. Our goal here is to show that even a naive approach to clustering can bring value, rather than to perform a systematic study of clustering heuristics, which we leave to future work. In Appendix A.9, we illustrate an application of the CR-ADP approach to a real world system with 200 locations.

n	CR-ADP	Myo.	No-R	3-RH-M	5-RH-M	7-RH-M	10-RH-M
20	4.04	4.73	9.40	4.51	4.43	4.38	4.31
30	2.91	3.17	6.10	3.70	3.66	3.62	3.59
40	3.35	3.96	8.15	4.60	4.53	4.04	4.05
50	3.77	4.16	8.48	4.33	4.19	4.14	-
60	3.80	4.20	9.01	4.36	4.24	-	-
70	3.77	4.13	8.84	4.31	4.18	-	-
80	3.34	3.83	7.93	3.96	3.87	-	-
90	3.83	4.16	8.72	4.33	-	-	-
100	3.41	3.90	7.69	3.89	-	-	-

Table 3 Summary of Results on Large-scale Instances ($\tilde{n} = 10$, $\rho = 0.95$)

n	CR-ADP	Myo.	No-R	3-RH-M	5-RH-M	7-RH-M	10-RH-M
20	21.50	26.09	61.98	24.83	24.01	23.65	23.50
30	14.55	17.04	51.10	20.12	19.36	19.21	18.83
40	16.18	20.63	52.46	22.42	21.97	21.59	21.53
50	18.88	21.81	59.24	23.22	22.54	22.16	-
60	19.97	22.17	60.70	23.37	23.17	-	-
70	18.56	21.81	59.90	22.89	22.63	-	-
80	15.67	20.24	55.75	21.23	21.17	-	-
90	19.00	21.85	59.86	23.16	-	-	-
100	16.65	20.48	53.06	20.76	-	-	-

Table 4 Summary of Results on Large-scale Instances ($\tilde{n} = 10$, $\rho = 0.99$)

8. Conclusion

In this paper, we consider the problem of optimal repositioning of inventory in a product rental network with multiple locations and where demand, rental periods, and return locations are stochastic. We show that the optimal policy is specified in terms of a region in the state space, inside of which it is optimal not to carry out any repositioning and outside of which it is optimal to reposition inventory. We also prove that when repositioning, it is always optimal to do so such that the system moves to a new state that is on the boundary of the no-repositioning region and provide a simple check for when a state is in the no-repositioning region. We then propose a provably convergent approximate dynamic programming algorithm, R-ADP, that builds a lower approximation of the convex value function by iteratively adding hyperplanes. We also propose a clustering-based approach that allows our ADP algorithm to solve large-scale problems. Numerical experiments on problems with up to 100 locations support the effectiveness of the algorithmic approach.

References

Balseiro SR, Brown DB, Chen C (2020) Dynamic pricing of relocating resources in large networks. *Management Science* DOI: <https://doi.org/10.1287/mnsc.2020.3735>.

- Banerjee S, Freund D, Lykouris T (2017) Pricing and optimization in shared vehicle systems: An approximation framework, arXiv preprint arXiv:1608.06819.
- Banerjee S, Kanoria Y, Qian P (2018) Dynamic assignment control of a closed queueing network under complete resource pooling. *arXiv e-prints* arXiv:1803.
- Bélanger V, Ruiz A, Soriano P (2019) Recent optimization models and trends in location, relocation, and dispatching of emergency medical vehicles. *European Journal of Operational Research* 272(1):1–23.
- Benjaafar S, Wu S, Liu H, Gunnarsson EB (2021) Dimensioning on-demand vehicle sharing systems. *Management Science* DOI: <https://doi.org/10.1287/mnsc.2021.3957>.
- Berman O (1981) Dynamic repositioning of indistinguishable service units on transportation networks. *Transportation Science* 15(2):115–136.
- Bertsekas DP, Tsitsiklis JN (1996) *Neuro-Dynamic Programming* (Belmont, MA: Athena Scientific).
- Bimpikis K, Candogan O, Saban D (2019) Spatial pricing in ride-sharing networks. *Operations Research* 67(3):744–769.
- Birge JR, Zhao G (2007) Successive linear approximation solution of infinite-horizon dynamic stochastic programs. *SIAM Journal on Optimization* 18(4):1165–1186.
- Boyd S, Boyd SP, Vandenberghe L (2004) *Convex optimization* (Cambridge university press).
- Braverman A, Dai JG, Liu X, Ying L (2019) Empty-car routing in ridesharing systems. *Operations Research* 67(5):1437–1452.
- Brown DB, Smith JE (2011) Dynamic portfolio optimization with transaction costs: Heuristics and dual bounds. *Management Science* 57(10):1752–1770.
- Bruglieri M, Colorni A, Luè A (2014) The vehicle relocation problem for the one-way electric vehicle sharing: an application to the milan case. *Procedia-Social and Behavioral Sciences* 111:18–27.
- Chen ZL, Powell WB (1999) Convergent cutting-plane and partial-sampling algorithm for multistage stochastic linear programs with recourse. *Journal of Optimization Theory and Applications* 102(3):497–524.
- Chung H, Freund D, Shmoys DB (2018) Bike Angels: An analysis of Citi Bike’s incentive program. *Proceedings of the 1st ACM SIGCAS Conference on Computing and Sustainable Societies* (ACM).
- Constantinides GM (1979) Multiperiod consumption and investment behavior with convex transactions costs. *Management Science* 25(11):1127–1137.
- De Farias DP, Van Roy B (2000) On the existence of fixed points for approximate value iteration and temporal-difference learning. *Journal of Optimization Theory and Applications* 105(3):589–608.
- Eberly JC, Van Mieghem JA (1997) Multi-factor dynamic investment under uncertainty. *Journal of economic theory* 75(2):345–387.
- Freund D, Henderson SG, Shmoys DB (2017) Minimizing multimodular functions and allocating capacity in bike-sharing systems. *International Conference on Integer Programming and Combinatorial Optimization*, 186–198 (Springer).

- Freund D, Henderson SG, Shmoys DB (2019) Bike sharing. Hu M, ed., *Sharing Economy: Making Supply Meet Demand* (Springer).
- Freund D, Norouzi-Fard A, Paul A, Henderson SG, Shmoys DB (2016) Data-driven rebalancing methods for bike-share systems. Technical report.
- George DK, Xia CH (2011) Fleet-sizing and service availability for a vehicle rental system via closed queueing networks. *European Journal of Operational Research* 211:198–207.
- Ghosh S, Varakantham P, Adulyasak Y, Jaillet P (2017) Dynamic repositioning to reduce lost demand in bike sharing systems. *Journal of Artificial Intelligence Research* 58:387–430.
- Girardeau P, Leclere V, Philpott AB (2014) On the convergence of decomposition methods for multistage stochastic convex programs. *Mathematics of Operations Research* 40(1):130–145.
- Godfrey GA, Powell WB (2001) An adaptive, distribution-free algorithm for the newsvendor problem with censored demands, with applications to inventory and distribution. *Management Science* 47(8):1101–1112.
- He L, Hu Z, Zhang M (2020) Robust repositioning for vehicle sharing. *Manufacturing & Service Operations Management* 22(2):241–256.
- He L, Mak HY, Rong Y (2019) Operations management of vehicle sharing systems. Hu M, ed., *Sharing Economy: Making Supply Meet Demand* (Springer).
- He L, Mak HY, Rong Y, Shen ZJM (2017) Service region design for urban electric vehicle sharing systems. *Manufacturing & Service Operations Management* 19(2):309–327.
- Higle JL, Sen S (1991) Stochastic decomposition: An algorithm for two-stage linear programs with recourse. *Mathematics of Operations Research* 16(3):650–669.
- Jian N, Freund D, Wiberg HM, Henderson SG (2016) Simulation optimization for a large-scale bike-sharing system. *Proceedings of the 2016 Winter Simulation Conference*, 602–613 (IEEE Press).
- Jiang N, Kulesza A, Singh S, Lewis R (2015) The dependence of effective planning horizon on model accuracy. *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, 1181–1189 (Citeseer).
- Kabra A, Belavina E, Girotra K (2020) Bike-share systems: Accessibility and availability. *Management Science* 66(9):3803–3824.
- Kaspi M, Raviv T, Tzur M (2017) Bike-sharing systems: User dissatisfaction in the presence of unusable bicycles. *IIEE Transactions* 49(2):144–158.
- Kushner HJ, Yin GG (2003) *Stochastic approximation and recursive algorithms and applications*, volume 35 (Springer).
- Lee CY, Meng Q (2015) *Handbook of Ocean Container Transport Logistics* (Springer).

- Leland HE (1999) Optimal portfolio management with transactions costs and capital gains taxes, working Paper RPF-290, IBER, UC Berkeley.
- Li J, Leung CS, Wu Y, Liu K (2007) Allocation of empty containers between multi-ports. *European Journal of Operational Research* 182(1):400–412.
- Li Y, Zheng Y, Yang Q (2018) Dynamic bike reposition: A spatio-temporal reinforcement learning approach. *Proceedings of the 24th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (ACM).
- Li Z, Tao F (2010) On determining optimal fleet size and vehicle transfer policy for a car rental company. *Computers & operations research* 37(2):341–350.
- Lin LJ (1992) Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning* 8(3-4):293–321.
- Linowsky K, Philpott AB (2005) On the convergence of sampling-based decomposition algorithms for multistage stochastic programs. *Journal of Optimization Theory and Applications* 125(2):349–366.
- Liu J, Sun L, Chen W, Xiong H (2016) Rebalancing bike sharing systems: A multi-source data smart optimization. *Proceedings of the 22nd ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 1005–1014 (ACM).
- Lu M, Chen Z, Shen S (2017) Optimizing the profitability and quality of service in carshare systems under demand uncertainty. *Manufacturing & Service Operations Management* 20(2):162–180.
- Ma H, Fang F, Parkes DC (2020) Spatio-temporal pricing for ridesharing platforms. *ACM SIGecom Exchanges* 18(2):53–57.
- Maxwell MS, Ni EC, Tong C, Henderson SG, Topaloglu H, Hunter SR (2014) A bound on the performance of an optimal ambulance redeployment policy. *Operations Research* 62(5):1014–1027.
- Maxwell MS, Restrepo M, Henderson SG, Topaloglu H (2010) Approximate dynamic programming for ambulance redeployment. *INFORMS Journal on Computing* 22(2):266–281.
- Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller M, Fidjeland AK, Ostrovski G, et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518(7540):529.
- Munos R, Szepesvári C (2008) Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research* 9(May):815–857.
- Muthuraman K, Kumar S (2006) Multidimensional portfolio optimization with proportional transaction costs. *Mathematical Finance: An International Journal of Mathematics, Statistics and Financial Economics* 16(2):301–335.
- Nair R, Miller-Hooks E (2011) Fleet management for vehicle sharing operations. *Transportation Science* 45(4):524–540.
- Nascimento JM, Powell WB (2009) An optimal approximate dynamic programming algorithm for the lagged asset acquisition problem. *Mathematics of Operations Research* 34(1):210–237.

- Nascimento JM, Powell WB (2010) Dynamic programming models and algorithms for the mutual fund cash balance problem. *Management Science* 56(5):801–815.
- Nyotta B, Bravo F, Feldman J (2019) Free rides in dockless, electric vehicle sharing systems Available at SSRN: <https://ssrn.com/abstract=3391937>.
- O’Mahony E, Shmoys DB (2015) Data analysis and optimization for (citi) bike sharing. *AAAI*, 687–694.
- Pereira MVF, Pinto LMVG (1991) Multi-stage stochastic optimization applied to energy planning. *Mathematical Programming* 52:359–375.
- Philpott AB, Guan Z (2008) On the convergence of stochastic dual dynamic programming and related methods. *Operations Research Letters* 36(4):450–455.
- Porteus E (1982) Conditions for characterizing the structure of optimal strategies in infinite-horizon dynamic programs. *Journal of Optimization Theory and Applications* 36(3):419–432.
- Porteus EL (1975) On the optimality of structured policies in countable stage decision processes. *Management Science* 22(2):148–157.
- Powell WB (2007) *Approximate Dynamic Programming: Solving the Curses of Dimensionality* (John Wiley & Sons), 2nd edition.
- Powell WB, Ruszczyński A, Topaloglu H (2004) Learning algorithms for separable approximations of discrete stochastic optimization problems. *Mathematics of Operations Research* 29(4):814–836.
- Puterman ML (1994) *Markov Decision Processes: Discrete Stochastic Dynamic Programming* (New York, NY, USA: John Wiley & Sons, Inc.), 1st edition.
- Raviv T, Kolka O (2013) Optimal inventory management of a bike-sharing station. *IIE Transactions* 45(10):1077–1093.
- Rockafellar RT (1970) *Convex Analysis* (Princeton, NJ, USA: Princeton University Press), 1st edition.
- Schuijbroek J, Hampshire RC, Van Hoesel WJ (2017) Inventory rebalancing and vehicle routing in bike sharing systems. *European Journal of Operational Research* 257(3):992–1004.
- Shapiro A (2011) Analysis of stochastic dual dynamic programming method. *European Journal of Operational Research* 209(1):63–72.
- Shu J, Chou MC, Liu Q, Teo CP, Wang IL (2013) Models for effective deployment and redistribution of bicycles within public bicycle-sharing systems. *Operations Research* 61(6):1346–1359.
- Shui C, Szeto W (2018) Dynamic green bike repositioning problem—a hybrid rolling horizon artificial bee colony algorithm approach. *Transportation Research Part D: Transport and Environment* 60:119–136.
- Song DP (2005) Optimal threshold control of empty vehicle redistribution in two depot service systems. *IEEE Transactions on Automatic Control* 50(1):87–90.
- Sundararajan A (2016) *The sharing economy: The end of employment and the rise of crowd-based capitalism* (Mit Press).

- Van Mieghem JA (2003) Commissioned paper: Capacity management, investment, and hedging: Review and recent developments. *Manufacturing & Service Operations Management* 5(4):269–302.
- Warrington J, Beuchat PN, Lygeros J (2019) Generalized dual dynamic programming for infinite horizon problems in continuous state and action spaces. *IEEE Transactions on Automatic Control* 64(12):5012–5023.
- Waserhole A, Jost V (2016) Pricing in vehicle sharing systems: Optimization in queuing networks with product forms. *EURO Journal on Transportation and Logistics* 5(3):293–320.
- Zhao L, Liu Z, Hu P (2020) Dynamic repositioning for vehicle sharing with setup costs. *Operations Research Letters* 48(6):792–797.
- Zipkin PH (2000) *Foundations of Inventory Management* (McGraw-Hill, New York).

Appendix A: Appendix

A.1. Empirical Justification of Assumption 1 and Estimation of P_t

In this section, we describe our analysis of a publicly available Car2Go dataset¹⁵, which contains GPS locations of Car2Go rental vehicles, at roughly 15-minute intervals, between June 2012 and December 2013 in Portland, Oregon. The goal is to understand whether Assumption 1 is realistic in practice. We use the following procedure to obtain estimates of the return fraction p_t specified in Assumption 1, calculated using data from the entire year of 2013.

1. We assume hourly rental periods (i.e., repositioning occurs every hour) and that the city is divided into n regions, $1, 2, \dots, n$.
2. For each vehicle, we extract all of its trips throughout the year under the assumption that trips end when the vehicle remains at the same GPS coordinates between successive measurements (15 minutes). This process yielded a total of 367,736 trips across 311 vehicles. Each trip is associated with a start time, return time, rental origin, and return location.
3. For each hour t , we directly estimate the return fraction for each region i by computing the set of rentals that are ongoing at hour t originating from region i (therefore, this includes any rental from i initiated at or before t and returned strictly after t). The return fraction from region i is then the fraction of these trips that are returned at hour $t + 1$ (to any return location).

Table 5 shows the results of this estimation procedure for $n = 9$. These $n = 9$ regions are constructed using a 3x3 grid, defined by latitude regions $(45.35^\circ, 45.50^\circ]$, $(45.50^\circ, 45.55^\circ]$, and $(45.55^\circ, 45.82^\circ]$, and the longitude regions $(-122.86^\circ, -122.65^\circ]$, $(-122.65^\circ, -122.60^\circ]$, and $(-122.60^\circ, -121.65^\circ]$.

Assumption 1 states that the return fraction p_t does not depend on the rental origin, but we compute these quantities as if the dependency on rental origin is allowed. However, it can be seen that the variation of the return fraction across regions is minor (i.e., the return fractions in each row are similar), thereby providing empirical evidence that Assumption 1 is not unreasonable.

¹⁵ <https://aaronparecki.com/car2go/>

We also note that the return fraction matrix \mathbf{P}_t can be estimated using the same procedure that is outlined above, assuming that GPS data is available. To obtain an estimate of $p_{t,ij}$, we would simply alter Step 3 above to calculate the fraction of ongoing trips from region i that are returned specifically to region j at hour $t + 1$.

Time Period Region	1	2	3	4	5	6	7	8	9
8am-9am	0.770	0.765	0.776	0.787	0.780	0.773	0.798	0.794	0.777
9am-10am	0.743	0.724	0.791	0.740	0.739	0.721	0.731	0.742	0.728
10am-11am	0.657	0.674	0.626	0.677	0.671	0.652	0.705	0.703	0.665
11am-12pm	0.644	0.681	0.642	0.671	0.681	0.660	0.668	0.697	0.671
12pm-1pm	0.693	0.659	0.664	0.661	0.684	0.676	0.655	0.705	0.655
1pm-2pm	0.642	0.664	0.714	0.642	0.682	0.636	0.650	0.696	0.686
2pm-3pm	0.697	0.650	0.611	0.655	0.688	0.668	0.684	0.693	0.630
3pm-4pm	0.656	0.647	0.651	0.680	0.671	0.644	0.660	0.668	0.650
4pm-5pm	0.667	0.654	0.616	0.664	0.676	0.651	0.633	0.632	0.630
5pm-6pm	0.664	0.696	0.644	0.700	0.727	0.690	0.687	0.706	0.698

Table 5 Estimates of Return Fraction p_t (see Assumption 1) from 2013 Car2Go GPS Data

A.2. Properties of the Repositioning Cost

In this subsection, We describe some properties of the repositioning cost $C(\cdot)$ defined in (1).

LEMMA 2. $C(\cdot)$ satisfies the following properties:

1. (Positive Homogeneity): $C(t\mathbf{z}) = tC(\mathbf{z})$ for all $t \geq 0$.
2. (Convexity): $C(\lambda\mathbf{z}_1 + (1 - \lambda)\mathbf{z}_2) \leq \lambda C(\mathbf{z}_1) + (1 - \lambda)C(\mathbf{z}_2)$ for all $\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{H}$ and $\lambda \in [0, 1]$.
3. (Sub-Additivity): $C(\mathbf{z}_1 + \mathbf{z}_2) \leq C(\mathbf{z}_1) + C(\mathbf{z}_2)$ for all $\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{H}$.
4. (Continuity): $C(\mathbf{z})$ is continuous in $\mathbf{z} \in \mathcal{H}$.

Proof. It is clear that the linear program (1) is bounded feasible. Therefore, an optimal solution to (1) exists and the strong duality holds. The dual linear program can be written as

$$C(\mathbf{z}) = \max_{\boldsymbol{\lambda}} \boldsymbol{\lambda}^T \mathbf{z} \quad \text{subject to } \lambda_j - \lambda_i \leq c_{ij} \quad \forall \quad i, j. \quad (23)$$

From (23), we have

$$C(t\mathbf{z}) = \max \{t\boldsymbol{\lambda}^T \mathbf{z} : \lambda_j - \lambda_i \leq c_{ij}, \forall i, j\} = t \max \{\boldsymbol{\lambda}^T \mathbf{z} : \lambda_j - \lambda_i \leq c_{ij}, \forall i, j\} = tC(\mathbf{z}).$$

Therefore, $C(\cdot)$ is positively homogeneous. As the pointwise supremum of a collection of convex and lower semicontinuous functions ($\boldsymbol{\lambda}^T \mathbf{z}$ for each $\boldsymbol{\lambda}$), $C(\cdot)$ is also convex and lower semicontinuous. It is well known that a convex function on a locally simplicial convex set is upper semicontinuous (Rockafellar (1970) Theorem 10.2). Therefore, as \mathcal{H} is a polyhedron, $C(\cdot)$ must be continuous. From the positive homogeneity and the convexity, we have

$$C(\mathbf{z}_1 + \mathbf{z}_2) = 2C\left(\frac{1}{2}\mathbf{z}_1 + \frac{1}{2}\mathbf{z}_2\right) \leq 2\left(\frac{1}{2}C(\mathbf{z}_1) + \frac{1}{2}C(\mathbf{z}_2)\right) = C(\mathbf{z}_1) + C(\mathbf{z}_2).$$

Therefore, $C(\cdot)$ is sub-additive.

Moreover, due to the triangle inequality, it is not optimal to simultaneously move inventory into and out of the same location. This property can be stated as follows.

LEMMA 3. *There exists an optimal solution \mathbf{w} to (1) such that*

$$\sum_{i=1}^n w_{ij} = z_j^+ \quad \text{and} \quad \sum_{k=1}^n w_{jk} = z_j^- \quad \text{for all } j = 1, \dots, n.$$

Proof. It is easy to see that an equivalent condition is $w_{i,j}w_{j,k} = 0$ for all i, j, k . To show this is true, suppose \mathbf{w} is an optimal solution and there exists i, j, k such that $w_{i,j}, w_{j,k} > 0$. If $i = k$, we can set at least one of $w_{i,j}$ and $w_{j,i}$ to 0 without violating the constraints. If $i \neq k$, we can set at least one of $w_{i,j}$ and $w_{j,k}$ to 0, and increase $w_{i,k}$ accordingly. In both cases, the resulting objective is at least as good. Repeating this for all i, k and j can enforce this condition for all i, k and j .

Lemma 3 leads to the following bound for the repositioning cost $C(\mathbf{z})$.

LEMMA 4.

$$\frac{c_{\min}}{2} \sum_{i=1}^n |z_i| \leq C(\mathbf{z}) \leq \frac{c_{\max}}{2} \sum_{i=1}^n |z_i|. \quad (24)$$

The proof follows from Lemma 3. There exists an optimal solution \mathbf{w} to (1) such that

$$\begin{aligned} C(\mathbf{z}) &= \sum_{i,j} c_{ij} w_{ij} = \frac{1}{2} \sum_j \sum_i c_{ij} w_{ij} + \frac{1}{2} \sum_j \sum_i c_{ji} w_{ji} \\ &\leq \frac{c_{\max}}{2} \sum_j z_j^+ + \frac{c_{\max}}{2} \sum_j z_j^- = \frac{c_{\max}}{2} \sum_j |z_j|. \end{aligned}$$

It is easy to see that, in (24), the equality holds if $c_{ij} = c_{\max}$ for all i, j . Therefore, the bound is tight. The lower bound follows the same logic.

A.3. Proof of Theorem 1

In this subsection, we provide a complete and self-contained proof for our main result, Theorem 1. We first provide some basic properties for directional derivatives. Recall the definition of the directional derivative:

$$u'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}, \boldsymbol{\eta}) = \lim_{t \downarrow 0} \frac{u(\mathbf{x} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) - u(\mathbf{x}, \boldsymbol{\gamma})}{t}.$$

LEMMA 5. *If $u(\mathbf{x}, \boldsymbol{\gamma})$ is jointly convex in \mathbf{x} and $\boldsymbol{\gamma}$, then $u'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}, \boldsymbol{\eta})$ satisfies the following properties:*

- (Positive Homogeneity) $u'(\mathbf{x}, \boldsymbol{\gamma}; t\mathbf{z}, t\boldsymbol{\eta}) = tu'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}, \boldsymbol{\eta})$
- (Sub-Additivity) $u'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}_1 + \mathbf{z}_2, \boldsymbol{\eta}_1 + \boldsymbol{\eta}_2) \leq u'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}_1, \boldsymbol{\eta}_1) + u'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{z}_2, \boldsymbol{\eta}_2)$.

Proof. These are well-known important properties for directional derivatives (see, for example, Rockafellar (1970)).

Now, we establish the fact that if $v_{t+1}(\cdot)$ is convex and additionally satisfies certain bounds on its directional derivatives, then for any realization of $\mathbf{d}_t, \mathbf{P}_t$, the function $U_t(\mathbf{y}_t, \boldsymbol{\gamma}_t, \mathbf{d}_t, \mathbf{P}_t)$ can be reformulated as the convex program (8).

PROPOSITION 3. *Suppose Assumptions 1 and 2 hold and*

1. $v_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1})$ is continuous and jointly convex in \mathbf{x}_{t+1} and $\boldsymbol{\gamma}_{t+1}$,

2. $v'_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}; \mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta}) \leq C(-\mathbf{z}) + \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) \eta_i$ for any feasible direction $(\mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$.

Then for all realizations (\mathbf{d}, \mathbf{p}) , $U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t)$ defined in (5) can be written as the following optimization problem:

$$\begin{aligned} U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t) = \min_{\mathbf{w}, \mathbf{x}_{t+1}, \gamma_{t+1}} \quad & \sum_{i=1}^n \beta_i (d_{t,i} - w_i) + \rho v_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}) \\ \text{s.t.} \quad & \mathbf{x}_{t+1} = \mathbf{y}_t - \mathbf{w} + \mathbf{P}_t^T (\gamma_t + \mathbf{w}); \\ & \gamma_{t+1} = (\gamma_t + \mathbf{w}) \circ (\mathbf{e} - \mathbf{P}_t \mathbf{e}); \\ & \mathbf{w} \leq \mathbf{y}_t; \\ & \mathbf{w} \leq \mathbf{d}_t. \end{aligned} \quad (25)$$

Moreover, $U_t(\mathbf{y}_t, \gamma_t, \mathbf{d}_t, \mathbf{P}_t)$ is continuous and jointly convex in \mathbf{y}_t and γ_t .

Proof. Let $(\mathbf{w}^*, \mathbf{x}_{t+1}^*, \gamma_{t+1}^*)$ be the optimal solution to problem (25). Clearly, if $\mathbf{w}^* = \min\{\mathbf{y}_t, \mathbf{d}_t\}$, the reformulation is correct. Now suppose that there exist an index i such that $\mathbf{w}^* < \min\{\mathbf{y}_t, \mathbf{d}_t\}$. Let $\hat{\mathbf{w}} = \mathbf{w}^* + \epsilon \mathbf{e}_i$ for a small $\epsilon > 0$. Then we have

$$\begin{aligned} \hat{\mathbf{x}}_{t+1} &= \mathbf{y}_t - (\mathbf{w}^* + \epsilon \mathbf{e}_i) + \mathbf{P}_t^T (\gamma_t + \mathbf{w}^* + \epsilon \mathbf{e}_i) = \mathbf{x}_{t+1}^* + \epsilon(-\mathbf{e}_i + \mathbf{P}_t^T \mathbf{e}_i) \\ \hat{\gamma}_{t+1} &= (\gamma_t + \mathbf{w}^* + \epsilon \mathbf{e}_i) \circ (\mathbf{e} - \mathbf{P}_t \mathbf{e}) = \gamma_{t+1}^* + \epsilon \mathbf{e}_i \left(1 - \sum_{j=1}^n P_{t,ij} \right) = \gamma_{t+1}^* + \epsilon \mathbf{e}_i (1 - p_t), \end{aligned}$$

where the last inequality is based on Assumption 1. Let

$$\hat{\mathbf{z}} = \mathbf{P}_t^T \mathbf{e}_i - \sum_{j=1}^n P_{t,ij} \mathbf{e}_i = \mathbf{P}_t^T \mathbf{e}_i - p_t \mathbf{e}_i, \quad \hat{\boldsymbol{\eta}} = \mathbf{e}_i (1 - p_t).$$

From the problem assumption, we have:

$$\begin{aligned} & v'_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}; -\mathbf{e}_i + \mathbf{P}_t^T \mathbf{e}_i, \mathbf{e}_i (1 - p_t)) \\ &= v'_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}; \hat{\mathbf{z}} - \hat{\boldsymbol{\eta}}, \hat{\boldsymbol{\eta}}) \\ &\leq C(-\hat{\mathbf{z}}) + \sum_{j=1}^n (\beta_j + \rho c_{\max} - c_{\min}) \hat{\eta}_j \\ &\leq \frac{c_{\max}}{2} \|\hat{\mathbf{z}}\|_1 + (\beta_i + \rho c_{\max} - c_{\min}) (1 - p_t) \quad (\text{by Lemma 4}) \\ &\leq \frac{c_{\max}}{2} \left(\|\mathbf{P}_t^T \mathbf{e}_i\|_1 + \left\| \sum_{j=1}^n P_{t,ij} \mathbf{e}_i \right\|_1 \right) + (\beta_i + \rho c_{\max} - c_{\min}) (1 - p_t) \quad (\text{by the triangle inequality}) \\ &= c_{\max} p_t + (\beta_i + \rho c_{\max} - c_{\min}) (1 - p_t). \end{aligned}$$

Let

$$f(\mathbf{w}) = \sum_{i=1}^n \beta_i (d_{t,i} - w_i) + \rho v_{t+1}(\mathbf{x}_{t+1}(\mathbf{w}), \gamma_{t+1}(\mathbf{w})),$$

where $\mathbf{x}_{t+1}(\mathbf{w}), \gamma_{t+1}(\mathbf{w})$ are defined from the first two equalities. Then

$$\begin{aligned} \frac{\partial f(\mathbf{w})}{\partial w_i} &= -\beta_i + \rho v'_{t+1}(\mathbf{x}_{t+1}, \gamma_{t+1}; -\mathbf{e}_i + \mathbf{P}_t^T \mathbf{e}_i, \mathbf{e}_i (1 - p_t)) \\ &\leq -\beta_i + \rho c_{\max} p_t + \rho (\beta_i + \rho c_{\max} - c_{\min}) (1 - p_t) \\ &\leq -\beta_i + \rho c_{\max} p_t + (\beta_i + \rho c_{\max} - c_{\min}) (1 - p_t) \quad (\text{since } \rho \leq 1) \\ &= \rho c_{\max} - c_{\min} - p_t (\beta_i - c_{\min}) \leq 0 \quad (\text{by Assumption 2}). \end{aligned}$$

Therefore, it is always optimal to increase w_i for all $\mathbf{y}_t, \boldsymbol{\gamma}_t, \mathbf{d}_t, \mathbf{P}_t$. This prove the first part of the proposition. The continuity of $u(\cdot)$ follows from the Dominated Convergence Theorem, as $U_t(\mathbf{y}_t, \boldsymbol{\gamma}_t, \mathbf{d}_t, \mathbf{P}_t) \leq \sum_i \beta_i d_i + \rho \|v\|_\infty$. For the convexity, since the constraint of the optimization problem is linear, thus the feasible region is convex. Also $v_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1})$ is jointly convex in \mathbf{x}_{t+1} and $\boldsymbol{\gamma}_{t+1}$, we have $U_t(\mathbf{y}_t, \boldsymbol{\gamma}_t, \mathbf{d}_t, \mathbf{P}_t)$ is jointly convex in \mathbf{y}_t and $\boldsymbol{\gamma}_t$ (see, e.g., section 3.2.5 of Boyd et al. (2004)).

From Proposition 3, to show the convexity of $u_t(\cdot)$, we require $v_{t+1}(\cdot)$ to be convex and satisfy the bounds on directional derivative as presented in item 2. If the convexity of $u_t(\cdot)$ could imply the convexity and the aforementioned bounds of directional derivation for $v_t(\cdot)$, then the induction is complete. Unfortunately, the convexity of $u_t(\cdot)$ does not imply that

$$v'_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1}; \mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta}) \leq C(-\mathbf{z}) + \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) \eta_i.$$

To overcome this difficulty, in Proposition 4, we assume a stronger condition on $v_{t+1}(\cdot)$ that implies that $u_t(\cdot)$ satisfies two types of bounds on partial derivative as stated in the Theorem 1. Then in Proposition 5, we show that if $u_t(\cdot)$ is convex and satisfies these two types of bounds, $v_t(\cdot)$ would be convex and satisfies the stronger condition assumed in Proposition 4. This would complete the induction step.

Before we present Proposition 4, we first show how to decompose the directional derivative of $U_{t, \mathbf{d}, \mathbf{P}}(\cdot, \cdot) \triangleq U_t(\cdot, \cdot, \mathbf{d}, \mathbf{P})$. Through this Lemma, we connect the directional derivatives of $u_t(\cdot)$ and those of $v_{t+1}(\cdot)$. We define some notation for indices sets. For any vector $\mathbf{y} \in \mathbb{R}^n$, we let $J^-(\mathbf{y}) = \{i | y_i < 0\}$, $J^0(\mathbf{y}) = \{i | y_i = 0\}$, $J^+(\mathbf{y}) = \{i | y_i > 0\}$, $J^{0+}(\mathbf{y}) = \{i | y_i \geq 0\}$ and $J^{0-}(\mathbf{y}) = \{i | y_i \leq 0\}$.

LEMMA 6. *For any realization (\mathbf{d}, \mathbf{P}) , we have*

$$U'_{t, \mathbf{d}, \mathbf{P}}(\mathbf{y}, \boldsymbol{\gamma}; \mathbf{z}, \boldsymbol{\eta}) = - \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} \beta_i z_i + \rho v'_{t+1}(\mathbf{x}, \boldsymbol{\zeta}; \mathbf{w}^+, \boldsymbol{\delta}^+), \quad (26)$$

where

$$\mathbf{w}_i^+ = \begin{cases} z_i + \iota_i + \sum_{j \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} z_j p_{ji} & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^+(\mathbf{z})), \\ \iota_i + \sum_{j \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} z_j p_{ji} & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^{0-}(\mathbf{z})), \end{cases}$$

with $\iota_i \triangleq \sum_{j=1}^n \eta_j p_{ji}$ and

$$\boldsymbol{\delta}_i^+ = \begin{cases} \eta_i (1 - \sum_{j=1}^n p_{ij}) & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^+(\mathbf{z})), \\ (\eta_i + z_i) (1 - \sum_{j=1}^n p_{ij}) & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^{0-}(\mathbf{z})), \end{cases}$$

and

$$\mathbf{x} = \tau_x(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}, \mathbf{P}); \quad \boldsymbol{\zeta} = \tau_\gamma(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}, \mathbf{P}).$$

Proof. Let

$$\vartheta_i = \sum_{j=1}^n \gamma_j p_{ji}, \quad \iota_i = \sum_{j=1}^n \eta_j p_{ji}$$

Note that

$$L(\mathbf{y}, \mathbf{d}) = \sum_{i \in J^-(\mathbf{y}-\mathbf{d})} \beta_i (d_i - y_i), \quad (27)$$

and let the next state, under (\mathbf{d}, \mathbf{P}) , be defined by $\mathbf{x}(\mathbf{y}, \boldsymbol{\gamma}), \boldsymbol{\zeta}(\mathbf{y}, \boldsymbol{\gamma})$, with components

$$x_i(\mathbf{y}, \boldsymbol{\gamma}) = \begin{cases} (y_i - d_i) + \vartheta_i + \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} d_j p_{ji} + \sum_{j \in J^0(\mathbf{y}-\mathbf{d})} y_j p_{ji} & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}), \\ \vartheta_i + \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} d_j p_{ji} + \sum_{j \in J^0(\mathbf{y}-\mathbf{d})} y_j p_{ji} & \text{for } i \in J^0(\mathbf{y}-\mathbf{d}), \end{cases} \quad (28)$$

$$\zeta_i(\mathbf{y}, \boldsymbol{\gamma}) = \begin{cases} (\gamma_i + d_i)(1 - \sum_{j=1}^n p_{ij}) & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}), \\ (\gamma_i + y_i)(1 - \sum_{j=1}^n p_{ij}) & \text{for } i \in J^0(\mathbf{y}-\mathbf{d}). \end{cases} \quad (29)$$

Choose t small enough so that the following hold:

$$\begin{aligned} J^-(\mathbf{y} + t\mathbf{z} - \mathbf{d}) &= J^-(\mathbf{y} - \mathbf{d}) \cup (J^0(\mathbf{y} - \mathbf{d}) \cap J^-(\mathbf{z})), \\ J^0(\mathbf{y} + t\mathbf{z} - \mathbf{d}) &= J^0(\mathbf{y} - \mathbf{d}) \cap J^0(\mathbf{z}), \\ J^+(\mathbf{y} + t\mathbf{z} - \mathbf{d}) &= J^+(\mathbf{y} - \mathbf{d}) \cup (J^0(\mathbf{y} - \mathbf{d}) \cap J^+(\mathbf{z})), \\ J^{0-}(\mathbf{y} + t\mathbf{z} - \mathbf{d}) &= J^-(\mathbf{y} - \mathbf{d}) \cup (J^0(\mathbf{y} - \mathbf{d}) \cap J^{0-}(\mathbf{z})), \end{aligned} \quad (30)$$

where the last equation follows directly from the first and second. For $\mathbf{y} + t\mathbf{z}$, we have, directly by (27), that

$$L(\mathbf{y} + t\mathbf{z}, \mathbf{d}) = \sum_{i \in J^-(\mathbf{y} + t\mathbf{z} - \mathbf{d})} \beta_i (d_i - y_i - tz_i),$$

and directly from (28) and (29),

$$\begin{aligned} &x_i(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) \\ &= \begin{cases} y_i + tz_i - d_i + \vartheta_i + t\iota_i + \sum_{j \in J^+(\mathbf{y} + t\mathbf{z} - \mathbf{d})} d_j p_{ji} & \text{for } i \in J^+(\mathbf{y} + t\mathbf{z} - \mathbf{d}), \\ + \sum_{j \in J^0(\mathbf{y} + t\mathbf{z} - \mathbf{d})} (y_j + tz_j) p_{ji} \\ \vartheta_i + t\iota_i + \sum_{j \in J^+(\mathbf{y} + t\mathbf{z} - \mathbf{d})} d_j p_{ji} & \text{for } i \in J^0(\mathbf{y} + t\mathbf{z} - \mathbf{d}), \\ + \sum_{j \in J^0(\mathbf{y} + t\mathbf{z} - \mathbf{d})} (y_j + tz_j) p_{ji} \end{cases} \end{aligned}$$

and

$$\zeta_i(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) = \begin{cases} (\gamma_i + t\eta_i + d_i)(1 - \sum_{j=1}^n p_{i,j}) & \text{for } i \in J^+(\mathbf{y} + t\mathbf{z} - \mathbf{d}), \\ (\gamma_i + t\eta_i + y_i + tz_i)(1 - \sum_{j=1}^n p_{i,j}) & \text{for } i \in J^0(\mathbf{y} + t\mathbf{z} - \mathbf{d}). \end{cases}$$

It follows by (30) that

$$L(\mathbf{y} + t\mathbf{z}, \mathbf{d}) - L(\mathbf{y}, \mathbf{d}) = -t \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} \beta_i z_i.$$

For the state update equations, we have

$$\begin{aligned} &x_i(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) - x_i(\mathbf{y}, \boldsymbol{\gamma}) \\ &= \begin{cases} tz_i + t\iota_i + \sum_{j \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} tz_j p_{ji} & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^+(\mathbf{z})), \\ t\iota_i + \sum_{j \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} tz_j p_{ji} & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^{0-}(\mathbf{z})), \end{cases} \end{aligned}$$

and

$$\begin{aligned} &\zeta_i(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) - \zeta_i(\mathbf{y}, \boldsymbol{\gamma}) \\ &= \begin{cases} t\eta_i(1 - \sum_{j=1}^n p_{ij}) & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^+(\mathbf{z})), \\ t(\eta_i + z_i)(1 - \sum_{j=1}^n p_{ij}) & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^{0-}(\mathbf{z})). \end{cases} \end{aligned}$$

Set

$$\boldsymbol{w}^+ = \frac{\mathbf{x}(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) - \mathbf{x}(\mathbf{y}, \boldsymbol{\gamma})}{t}, \quad \boldsymbol{\delta}^+ = \frac{\boldsymbol{\zeta}(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}) - \boldsymbol{\zeta}(\mathbf{y}, \boldsymbol{\gamma})}{t}.$$

Then

$$\lim_{t \rightarrow 0} \frac{v_{t+1}(\mathbf{x}(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta}), \boldsymbol{\zeta}(\mathbf{y} + t\mathbf{z}, \boldsymbol{\gamma} + t\boldsymbol{\eta})) - v_{t+1}(\mathbf{x}(\mathbf{y}, \boldsymbol{\gamma}), \boldsymbol{\zeta}(\mathbf{y}, \boldsymbol{\gamma}))}{t} = v'_{t+1}(\mathbf{x}, \boldsymbol{\zeta}; \mathbf{w}^+, \boldsymbol{\delta}^+).$$

It follows that

$$U'_{t,\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; \mathbf{z}, \boldsymbol{\eta}) = - \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup (J^0(\mathbf{y}-\mathbf{d}) \cap J^-(\mathbf{z}))} \beta_i z_i + \rho v'_{t+1}(\mathbf{x}, \boldsymbol{\zeta}; \mathbf{w}^+, \boldsymbol{\delta}^+). \quad (31)$$

This concludes the proof.

Now we are ready to present Proposition 4, which shows that given the convexity of $v_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1})$ and certain bounds on its directional derivatives, the function $u_t(\mathbf{y}_t, \boldsymbol{\gamma}_t)$ satisfies two types of bounds on its directional derivatives.

PROPOSITION 4. *Suppose Assumptions 1 and 2 hold and*

1. $v_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1})$ is continuous and jointly convex in \mathbf{x}_{t+1} and $\boldsymbol{\gamma}_{t+1}$,
2. $v'_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1}; \mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta}) \leq C(-\mathbf{z}) + \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) \eta_i$ for any feasible direction $(\mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$,
3. $v'_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1}; \mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta}) \leq C(-\mathbf{z}) + \sum_{i=1}^n \beta_i \eta_i$ for any feasible direction $(\mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$,
4. $v'_{t+1}(\mathbf{x}_{t+1}, \boldsymbol{\gamma}_{t+1}; \mathbf{0}, \mathbf{z}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |z_i|$ for any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = \mathbf{0}$,

then we have:

1. $|u'_t(\mathbf{y}_t, \boldsymbol{\gamma}_t; \mp \boldsymbol{\xi}, \pm \boldsymbol{\xi})| \leq \sum_{i=1}^n \beta_i \xi_i$ for all $(\mathbf{y}_t, \boldsymbol{\gamma}_t) \in \Delta$ and any feasible direction $(\mp \boldsymbol{\xi}, \pm \boldsymbol{\xi})$ with $\boldsymbol{\xi} \geq \mathbf{0}$;
2. $u'_t(\mathbf{y}_t, \boldsymbol{\gamma}_t; \mathbf{0}, \mathbf{z}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |z_i|$ for all $(\mathbf{y}_t, \boldsymbol{\gamma}_t) \in \Delta$ and any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = \mathbf{0}$.

Proof. We omit the subscript t to reduce notation. To show the first inequality, we start with showing that $u'(\mathbf{y}, \boldsymbol{\gamma}; -\boldsymbol{\xi}, \boldsymbol{\xi}) \leq \sum_i \beta_i \xi_i$. From Lemma 6, noting that $-\boldsymbol{\xi} \leq \mathbf{0}$, we have

$$U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; -\boldsymbol{\xi}, \boldsymbol{\xi}) = \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i \xi_i + \rho v'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{w} - \boldsymbol{\delta}, \boldsymbol{\delta}),$$

where

$$w_i = \begin{cases} -\xi_i \sum_{j=1}^n p_{ij} + \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} \xi_j p_{ji} & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}), \\ \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} \xi_j p_{ji} & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d}), \end{cases}$$

and

$$\delta_i = \begin{cases} \xi_i (1 - \sum_{j=1}^n p_{ij}) = \xi_i (1 - \rho) & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}), \\ 0 & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d}). \end{cases}$$

Note that $\mathbf{e}^T \mathbf{w} = 0$ and so it is clear that $(\mathbf{w} - \boldsymbol{\delta}, \boldsymbol{\delta})$ is a feasible direction at $(\mathbf{x}, \boldsymbol{\zeta})$. It follows that

$$\begin{aligned} & U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; -\boldsymbol{\xi}, \boldsymbol{\xi}) \\ &= \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i \xi_i + \rho v'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{w} - \boldsymbol{\delta}, \boldsymbol{\delta}) \\ &\leq \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \rho C(-\mathbf{w}) + \rho \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) |\delta_i| \quad (\text{by assumption}) \\ &\leq \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \frac{\rho c_{\max}}{2} \sum_{i=1}^n |w_i| + \rho \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) |\delta_i| \quad (\text{by Lemma 4}) \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \frac{2\rho c_{\max}}{2} \sum_{i=1}^n \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} |\xi_j| p_{ji} \\
&\quad + \rho \sum_{i \in J^+(\mathbf{y}-\mathbf{d})} (\beta_i + \rho c_{\max} - c_{\min}) |\xi_i| (1-p) \quad (\text{by the triangle inequality}) \\
&= \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} (p\rho c_{\max} + \rho(\beta_i + \rho c_{\max} - c_{\min})(1-p)) |\xi_j| \\
&\leq \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} (p\rho c_{\max} + (\beta_i + \rho c_{\max} - c_{\min})(1-p)) |\xi_j| \\
&\leq \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \sum_{j \in J^+(\mathbf{y}-\mathbf{d})} (\beta_i + (\rho c_{\max} - c_{\min}) - p(\beta_i - c_{\min})) |\xi_j| \\
&\leq \sum_{i \in J^-(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| + \sum_{i \in J^+(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| \quad (\text{by Assumption 2}) \\
&= \sum_{i=1}^n \beta_i |\xi_i|
\end{aligned}$$

So, $U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \gamma; -\xi, \xi) \leq \sum_{i=1}^n \beta_i |\xi_i| = \sum_{i=1}^n \beta_i \xi_i$ holds for each $(\mathbf{y}, \gamma) \in \Delta$ and $(z - \boldsymbol{\eta}, \boldsymbol{\eta})$ feasible. It follows that $u'(\mathbf{y}, \gamma; -\xi, \xi) = \int U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \gamma; -\xi, \xi) d\mu \leq \sum_i \beta_i \xi_i$. From Lemma 3, $u(\mathbf{y}, \gamma)$ is convex, thus

$$u'(\mathbf{y}, \gamma; \xi, -\xi) \geq -u'(\mathbf{y}, \gamma; -\xi, \xi) \geq -\sum_i \beta_i \xi_i.$$

Now we show that $u'(\mathbf{y}, \gamma; \xi, -\xi) \leq \sum_{i=1}^n \beta_i \xi_i$ for all $(\mathbf{y}, \gamma) \in \Delta$ for all feasible direction $(\xi, -\xi)$ with $\xi \geq 0$. From the previous analysis, we have

$$U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \gamma; \xi, -\xi) = -\sum_{i \in J^-(\mathbf{y}-\mathbf{d})} \beta_i \xi_i + \rho v'(\mathbf{x}, \zeta; \mathbf{w} + \boldsymbol{\delta}, -\boldsymbol{\delta})$$

where

$$\begin{aligned}
w_i &= \begin{cases} \xi_i \sum_{j=1}^n p_{ij} - \sum_{j \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \xi_j p_{ji} & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d}), \\ -\sum_{j \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \xi_j p_{ji} & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}), \end{cases} \\
\delta_i &= \begin{cases} \xi_i (1 - \sum_{j=1}^n p_{ij}) = \xi_i (1-p) & \text{for } i \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d}), \\ 0 & \text{for } i \in J^-(\mathbf{y}-\mathbf{d}). \end{cases}
\end{aligned}$$

Clear $(\mathbf{w} + \boldsymbol{\delta}, -\boldsymbol{\delta})$ is a feasible direction at (\mathbf{x}, ζ) . It follows that

$$\begin{aligned}
&U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \gamma; \xi, -\xi) \\
&= -\beta_i \sum_{i \in J^-(\mathbf{y}-\mathbf{d})} \xi_i + \rho v'(\mathbf{x}, \gamma; \mathbf{w} + \boldsymbol{\delta}, -\boldsymbol{\delta}) \\
&\leq \frac{\rho c_{\max}}{2} \sum_{i=1}^n |w_i| + \rho \sum_{i=1}^n \beta_i |\delta_i| \quad (\text{from the third inequality of the assumption and } \xi \geq 0) \\
&\leq \frac{1}{2} \sum_{i=1}^n \beta_i |w_i| + \rho \sum_{i=1}^n \beta_i |\delta_i| \quad (\text{since } \rho c_{\max} \leq \beta_i, \rho \leq 1) \\
&\leq \frac{\rho}{2} 2 \sum_{i=1}^n \sum_{j \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_j| p_{ji} + \rho \sum_{i \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| \left(1 - \sum_{j=1}^n p_{ij}\right) \quad (\text{by the triangle inequality}) \\
&\leq \rho \sum_{i \in J^+(\mathbf{y}-\mathbf{d}) \cup J^0(\mathbf{y}-\mathbf{d})} \beta_i |\xi_i| \leq \sum_{i=1}^n \beta_i |\xi_i|.
\end{aligned}$$

So, $U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; \boldsymbol{\xi}, -\boldsymbol{\xi}) \leq \sum_{i=1}^n \beta_i |\xi_i| = \sum_{i=1}^n \beta_i \xi_i$ holds for each $(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta$ and $(\mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta})$ feasible. It follows that $u'(\mathbf{y}, \boldsymbol{\gamma}; \boldsymbol{\xi}, -\boldsymbol{\xi}) = \int U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; \boldsymbol{\xi}, -\boldsymbol{\xi}) d\mu \leq \sum_i \beta_i \xi_i$. From Lemma 3, $u(\mathbf{y}, \boldsymbol{\gamma})$ is convex, thus

$$u'(\mathbf{y}, \boldsymbol{\gamma}; -\boldsymbol{\xi}, \boldsymbol{\xi}) \geq -u'(\mathbf{y}, \boldsymbol{\gamma}; \boldsymbol{\xi}, -\boldsymbol{\xi}) \geq -\sum_i \beta_i \xi_i.$$

Summing up all the above, we have shown the first inequality.

Now we show the second inequality. From Lemma 6, we have

$$U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; \mathbf{0}, \mathbf{z}) = \rho v'(\mathbf{x}, \boldsymbol{\gamma}; \boldsymbol{\iota}, \mathbf{z}(1-p))$$

where $\iota_i = \sum_{j=1}^n z_j p_{ji} \forall i$ and $p = \sum_{j=1}^n p_{ij}$ (Assumption 1). Therefore,

$$\sum_{i=1}^n \iota_i = \sum_{j=1}^n \sum_{i=1}^n z_j p_{ji} = \sum_{j=1}^n p z_j = 0,$$

and $\mathbf{e}^T \mathbf{z}(1-p) = 0$. It follows that

$$\begin{aligned} U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; \mathbf{0}, \mathbf{z}) &= \rho v'(\mathbf{x}, \boldsymbol{\gamma}; \boldsymbol{\iota}, \mathbf{z}(1-p)) \\ &\leq \rho v'(\mathbf{x}, \boldsymbol{\gamma}; \boldsymbol{\iota}, \mathbf{0}) + \rho v'(\mathbf{x}, \boldsymbol{\gamma}; \mathbf{0}, \mathbf{z}(1-p)) \text{ (by subadditivity, Lemma 5)} \\ &\leq \rho C(-\boldsymbol{\iota}) + \frac{\rho^2 c_{\max}}{2} \sum_{i=1}^n |z_i| (1-p) \text{ (by assumption)} \\ &\leq \frac{\rho c_{\max}}{2} \sum_{i=1}^n |\iota_i| + \frac{\rho^2 c_{\max}}{2} \sum_{i=1}^n |z_i| (1-p) \text{ (by Lemma 4)} \\ &\leq \frac{\rho c_{\max}}{2} \sum_{i=1}^n \sum_{j=1}^n |z_j| p_{ji} + \frac{\rho^2 c_{\max}}{2} \sum_{i=1}^n |z_i| (1-p) \text{ (by the triangle inequality)} \\ &\leq \frac{\rho c_{\max}}{2} \sum_{i=1}^n |z_i|. \text{ (by the subadditivity)} \end{aligned}$$

So, $U'_{\mathbf{d},\mathbf{p}}(\mathbf{y}, \boldsymbol{\gamma}; \mathbf{0}, \mathbf{z}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |z_i|$ holds for all realizations (\mathbf{d}, \mathbf{p}) . It follows that the integral also satisfies the condition.

The upcoming result, Proposition 5, assists with the eventual induction over t by stating that if $u_t(\mathbf{y}_t, \boldsymbol{\gamma}_t)$ is convex and certain bounds on its directional derivatives are satisfied, then $v_t(\mathbf{x}_t, \boldsymbol{\gamma}_t)$ not only is convex, but also satisfies the bounds on the directional derivatives required by Proposition 4. Before continuing, we introduce a technical lemma that carefully analyzes the optimal repositioning plan. The main difficulty in showing Proposition 5 is that there exist some situations where $x_{ti} > 0$, but $\mathbf{y}_{ti}^* = \mathbf{0}$ where \mathbf{y}_t^* is the optimal solution to $v_t(\mathbf{x}_t, \boldsymbol{\gamma}_t)$. In this case, if we move along the direction \mathbf{z} with $z_i < 0$ starting from \mathbf{x}_t , we are not able to move along this direction if we start from position \mathbf{y}^* . Therefore, we would need to find another direction such that it is feasible for position \mathbf{y}^* , but the cost of repositioning can be controlled. This is the main purpose of Lemma 7. This lemma plays a crucial role in the proof of Proposition 5.

LEMMA 7. *Suppose $\mathbf{x}, \mathbf{y} \geq 0$ and $\mathbf{e}^T \mathbf{x} = \mathbf{e}^T \mathbf{y}$. Suppose $\mathbf{x} - \boldsymbol{\eta} \geq 0$. Then there exists another vector $\boldsymbol{\xi}$ such that: 1) $\mathbf{y} - \boldsymbol{\xi} \geq 0$, 2) $\mathbf{e}^T \boldsymbol{\xi} = \mathbf{e}^T \boldsymbol{\eta}$, 3) $\sum_{j=1}^n |\xi_j - \eta_j| \leq 2 \sum_{j=1}^n |\eta_j|$, 4) $C(\mathbf{y} - \boldsymbol{\xi} - \mathbf{x} + \boldsymbol{\eta}) = C(\mathbf{y} - \mathbf{x}) - C(\boldsymbol{\xi} - \boldsymbol{\eta})$. Moreover, if $\boldsymbol{\eta} \geq 0$, we have $\boldsymbol{\xi} \geq 0$ and*

$$\sum_{j=1}^n \xi_j + \frac{1}{2} \sum_{j=1}^n |\xi_j - \eta_j| \leq \sum_{j=1}^n \eta_j \tag{32}$$

Proof. Let $\mathcal{I} = \{i \mid y_i - \eta_i < 0\}$. If $\mathcal{I} = \emptyset$, we have $\mathbf{y} - \boldsymbol{\eta} \geq 0$ and we just find $\boldsymbol{\xi} = \boldsymbol{\eta}$. Thus in the followings, we assume that $\mathcal{I} \neq \emptyset$. We let $\mathcal{J}^- = \{i \mid y_i - x_i \leq 0\}$ and $\mathcal{J}^+ = \{i \mid y_i - x_i > 0\}$. Note that $y_i - x_i < 0 \forall i \in \mathcal{I}$, so $\mathcal{I} \subseteq \mathcal{J}^-$.

Lemma 3 suggests that there exists $\mathbf{w} \geq 0$ such that $\mathbf{c} \cdot \mathbf{w} = C(\mathbf{y} - \mathbf{x})$ and

$$y_j - x_j = \begin{cases} -\sum_{i \in \mathcal{J}^+} w_{ji} = -\sum_i w_{ji} & \text{for } j \in \mathcal{J}^-; \\ \sum_{i \in \mathcal{J}^-} w_{ij} = \sum_i w_{ij} & \text{for } j \in \mathcal{J}^+. \end{cases}$$

The interpretation of the first case above is that locations j with more inventory than the target level y_j transfer inventory to locations that are below the target level. The second case is interpreted in an analogous way. Let $\boldsymbol{\xi}$ be such that

$$\xi_j = \begin{cases} y_j & \text{for } j \in \mathcal{I}; \\ \eta_j & \text{for } j \in \mathcal{J}^- \setminus \mathcal{I}; \\ \sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} + \eta_j & \text{for } j \in \mathcal{J}^+. \end{cases}$$

We claim that $\boldsymbol{\xi}$ satisfies the desired properties. We first verify that $y_j - \xi_j \geq 0$. This is clearly true if $j \in \mathcal{J}^-$ from the construction of $\boldsymbol{\xi}$. Now let $j \in \mathcal{J}^+$, then we have

$$y_j - \xi_j = y_j - \sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} - \eta_j \geq y_j - \sum_{i \in \mathcal{I}} w_{ij} - \eta_j \geq y_j - \sum_i w_{ij} - \eta_j = x_j - \eta_j \geq 0,$$

where the first inequality follows from $x_i \geq \eta_i > y_i$ for $i \in \mathcal{I}$. Therefore, $\mathbf{y} - \boldsymbol{\xi} \geq 0$ and part (1) is complete.

Also, using $x_j - y_j = \sum_{i \in \mathcal{J}^+} w_{ij}$, we have:

$$\begin{aligned} \sum_j \xi_j &= \sum_{j \in \mathcal{I}} y_j + \sum_{j \in \mathcal{J}^- \setminus \mathcal{I}} \eta_j + \sum_{j \in \mathcal{J}^+} \left(\sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} + \eta_j \right) \\ &= \sum_{j \in \mathcal{I}} y_j + \sum_{j \in (\mathcal{J}^- \setminus \mathcal{I}) \cup \mathcal{J}^+} \eta_j + \sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} (x_i - y_i) = \sum_i \eta_i, \end{aligned}$$

verifying part (2). Note that:

$$\xi_j - \eta_j = \begin{cases} y_j - \eta_j & \text{for } j \in \mathcal{I}; \\ 0 & \text{for } j \in \mathcal{J}^- \setminus \mathcal{I}; \\ \sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} & \text{for } j \in \mathcal{J}^+. \end{cases}$$

To show part (3), we have

$$\sum_{j=1}^n |\xi_j - \eta_j| = 2 \sum_{j \in \mathcal{I}} (\eta_j - y_j) \leq 2 \sum_{j=1}^n |\eta_j|.$$

To show that $C(\mathbf{y} - \boldsymbol{\xi} - \mathbf{x} + \boldsymbol{\eta}) \leq C(\mathbf{y} - \mathbf{x}) - C(\boldsymbol{\xi} - \boldsymbol{\eta})$, let $\bar{\mathbf{w}}$ be such that

$$\bar{w}_{ij} = \begin{cases} (1 - \frac{\eta_i - y_i}{x_i - y_i}) w_{ij} & \text{for } i \in \mathcal{I}; \\ w_{ij} & \text{for } i \notin \mathcal{I}. \end{cases}$$

Then for all $j \in \mathcal{J}^+$, we have

$$\begin{aligned} & y_j - \xi_j - (x_j - \eta_j) - \sum_i \bar{w}_{ij} + \sum_i \bar{w}_{ji} \\ &= (y_j - x_j) - \xi_j + \eta_j - \sum_{i \in \mathcal{I}} \bar{w}_{ij} - \sum_{i \notin \mathcal{I}} \bar{w}_{ij} \\ &= \sum_i w_{ij} - \sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} - \sum_i w_{ij} + \sum_{i \in \mathcal{I}} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} = 0, \end{aligned}$$

where we used $\sum_i \bar{w}_{ji} = 0$. Similarly, for all $j \in \mathcal{I}$, we have

$$\begin{aligned} & x_j - \eta_j - y_j + \xi_j - \sum_i \bar{w}_{ji} + \sum_i \bar{w}_{ij} \\ &= x_j - \eta_j - \sum_i \left(1 - \frac{\eta_j - y_j}{x_j - y_j}\right) w_{ji} \\ &= x_j - \sum_i w_{ji} - \eta_j - \frac{\eta_j - y_j}{x_j - y_j} (y_j - x_j) = 0 \end{aligned}$$

and for all $j \in \mathcal{J}^- \setminus \mathcal{I}$, we have

$$x_j - \eta_j - y_j + \xi_j - \sum_i \bar{w}_{ji} + \sum_i \bar{w}_{ij} = x_j - y_j - \sum_i w_{ji} = 0.$$

Therefore, we have shown that $\bar{\mathbf{w}}$ is a feasible solution to the optimization problem for $C(\cdot)$ defined in (1) at $\mathbf{y} - \boldsymbol{\xi} - \mathbf{x} + \boldsymbol{\eta}$. Thus, we have:

$$\begin{aligned} C(\mathbf{y} - \boldsymbol{\xi} - \mathbf{x} + \boldsymbol{\eta}) &\leq \sum_{i,j} c_{ij} \bar{w}_{ij} \\ &= \sum_{i,j} c_{ij} w_{ij} - \sum_{i \in \mathcal{I}} \sum_j c_{ij} \left(\frac{\eta_i - y_i}{x_i - y_i} w_{ij} \right) \\ &= C(\mathbf{y} - \mathbf{x}) - \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}^+} c_{ij} \left(\frac{\eta_i - y_i}{x_i - y_i} w_{ij} \right). \end{aligned}$$

If we define $\hat{\mathbf{w}}$ as

$$\hat{w}_{ij} = \begin{cases} \frac{\eta_i - y_i}{x_i - y_i} w_{ij} & \text{for } i \in \mathcal{I}, j \in \mathcal{J}^+; \\ 0 & \text{otherwise,} \end{cases}$$

then we can check that $\hat{\mathbf{w}}$ is a feasible solution to (1) at $\boldsymbol{\xi} - \boldsymbol{\eta}$. Therefore, we have

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}^+} c_{ij} \left(\frac{\eta_i - y_i}{x_i - y_i} w_{ij} \right) = \mathbf{c} \cdot \hat{\mathbf{w}} \geq C(\boldsymbol{\xi} - \boldsymbol{\eta}).$$

Finally, we have

$$C(\mathbf{y} - \boldsymbol{\xi} - \mathbf{x} + \boldsymbol{\eta}) \leq C(\mathbf{y} - \mathbf{x}) - \sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{J}^+} c_{ij} \left(\frac{\eta_i - y_i}{x_i - y_i} w_{ij} \right) \leq C(\mathbf{y} - \mathbf{x}) - C(\boldsymbol{\xi} - \boldsymbol{\eta}).$$

Together with the subadditivity of $C(\cdot)$, we conclude that

$$C(\mathbf{y} - \boldsymbol{\xi} - \mathbf{x} + \boldsymbol{\eta}) = C(\mathbf{y} - \mathbf{x}) - C(\boldsymbol{\xi} - \boldsymbol{\eta}),$$

proving part (4). The last claim follows directly from the construction of $\boldsymbol{\xi}$ and parts (2) and (3).

Now we are ready to prove Proposition 5.

PROPOSITION 5. *Suppose $u_t(\cdot)$ is convex and continuous in Δ and*

1. $u'_t(\mathbf{y}_t, \boldsymbol{\gamma}_t; \pm \boldsymbol{\eta}, \mp \boldsymbol{\eta}) \leq \sum_{i=1}^n \beta_i \eta_i$ for all $(\mathbf{y}_t, \boldsymbol{\gamma}_t) \in \Delta$ and for any feasible direction $(\pm \boldsymbol{\eta}, \mp \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$;
2. $u'_t(\mathbf{y}_t, \boldsymbol{\gamma}_t; \mathbf{0}, \mathbf{z}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |z_i|$ for all $(\mathbf{y}_t, \boldsymbol{\gamma}_t) \in \Delta$ and for any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = 0$.

Then, the value function $v_t(\cdot)$ is convex and continuous in Δ with $\Omega_v(\gamma) = \Delta_{n-1}(I)$ for $\gamma \in S$. For each $(\mathbf{x}_t, \gamma_t) \in \Delta$ and $\boldsymbol{\eta} \geq \mathbf{0}$, the directional derivatives satisfy

$$v'_t(\mathbf{x}_t, \gamma_t; \mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta}) \leq C(-\mathbf{z}) + \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) \eta_i \quad (33)$$

for any feasible direction $(\mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta})$ and

$$v'_t(\mathbf{x}_t, \gamma_t; \mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta}) \leq C(-\mathbf{z}) + \sum_{i=1}^n \beta_i \eta_i \quad (34)$$

for any feasible direction $(\mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta})$. In addition,

$$v'_t(\mathbf{x}_t, \gamma_t; \mathbf{0}, \mathbf{z}) \leq \frac{\rho c_{\max}}{2} \sum_{i=1}^n |z_i| \quad (35)$$

for any feasible direction $(\mathbf{0}, \mathbf{z})$ with $\mathbf{e}^T \mathbf{z} = 0$.

Proof. We omit the subscript t throughout the proof to reduce notation. To show that $v(\cdot)$ is convex, suppose \mathbf{y}_1 and \mathbf{y}_2 are optimal solutions of (16) for (\mathbf{x}_1, γ_1) and (\mathbf{x}_2, γ_2) , respectively. Then, $\lambda \mathbf{y}_1 + (1 - \lambda) \mathbf{y}_2 \in \Delta_{n-1}(\lambda \mathbf{e}^T \mathbf{x}_1 + (1 - \lambda) \mathbf{e}^T \mathbf{x}_2)$ and thus

$$\begin{aligned} v(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2, \lambda \gamma_1 + (1 - \lambda) \gamma_2) &\leq u(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2, \lambda \gamma_1 + (1 - \lambda) \gamma_2) \\ &\quad + C(\lambda(\mathbf{y}_1 - \mathbf{x}_1) + (1 - \lambda)(\mathbf{y}_2 - \mathbf{x}_2)) \\ &\leq \lambda v(\mathbf{x}_1, \gamma_1) + (1 - \lambda) v(\mathbf{x}_2, \gamma_2), \end{aligned}$$

by convexity of $u(\cdot)$ and Lemma 2. Continuity follows from Berge's Maximum Theorem, as the set-valued map $\mathbf{x} \mapsto \Delta_{n-1}(I)$ is continuous. To show the result in (33), suppose $(\mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta})$ is a feasible direction. Let \mathbf{y}^* be an optimal solution to equation (16) at (\mathbf{x}, γ) . Therefore,

$$v(\mathbf{x}, \gamma) = \min_{\mathbf{y} \in \Delta_{n-1}(\mathbf{e}^T \mathbf{x})} C(\mathbf{y} - \mathbf{x}) + u(\mathbf{y}, \gamma) = C(\mathbf{y}^* - \mathbf{x}) + u(\mathbf{y}^*, \gamma).$$

Let $t > 0$ be small enough such that $\mathbf{x} + t(\mathbf{z} - \boldsymbol{\eta}) \geq \mathbf{0}$. According to Lemma 7, there exists a vector $\boldsymbol{\xi} \geq \mathbf{0}$ such that for small enough t : 1) $\mathbf{y}^* - t\boldsymbol{\xi} \geq \mathbf{0}$, 2) $\mathbf{e}^T \boldsymbol{\xi} = \mathbf{e}^T \boldsymbol{\eta}$, 3) 4) $C(\mathbf{y}^* - t\boldsymbol{\xi} - \mathbf{x} - t\mathbf{z} + t\boldsymbol{\eta}) = C(\mathbf{y}^* - \mathbf{x} - t\mathbf{z}) - tC(\boldsymbol{\xi} - \boldsymbol{\eta})$. Therefore, $\mathbf{y}^* - t\boldsymbol{\xi}$ is a feasible solution to equation (16) at $(\mathbf{x} + t\mathbf{z} - t\boldsymbol{\eta}, \gamma + t\boldsymbol{\eta})$ and thus we have

$$\begin{aligned} &\frac{v(\mathbf{x} + t\mathbf{z} - t\boldsymbol{\eta}, \gamma + t\boldsymbol{\eta}) - v(\mathbf{x}, \gamma)}{t} \\ &\leq \frac{u(\mathbf{y}^* - t\boldsymbol{\xi}, \gamma + t\boldsymbol{\eta}) - u(\mathbf{y}^* - t\boldsymbol{\xi}, \gamma + t\boldsymbol{\xi}) + u(\mathbf{y}^* - t\boldsymbol{\xi}, \gamma + t\boldsymbol{\xi}) - u(\mathbf{y}^*, \gamma)}{t} \\ &\quad + \frac{C(\mathbf{y}^* - t\boldsymbol{\xi} - \mathbf{x} - t\mathbf{z} + t\boldsymbol{\eta}) - C(\mathbf{y}^* - \mathbf{x} - t\mathbf{z}) + C(\mathbf{y}^* - \mathbf{x} - t\mathbf{z}) - C(\mathbf{y}^* - \mathbf{x})}{t} \\ &\leq \frac{u(\mathbf{y}^* - t\boldsymbol{\xi}, \gamma + t\boldsymbol{\eta}) - u(\mathbf{y}^* - t\boldsymbol{\xi}, \gamma + t\boldsymbol{\xi}) + u(\mathbf{y}^* - t\boldsymbol{\xi}, \gamma + t\boldsymbol{\xi}) - u(\mathbf{y}^*, \gamma)}{t} - C(\boldsymbol{\xi} - \boldsymbol{\eta}) + C(-\mathbf{z}) \end{aligned}$$

Adding and subtracting $u(\mathbf{y}^*, \gamma)$ in the numerator of the first term and then taking limits on both sides, we get

$$\begin{aligned} v'(\mathbf{x}, \gamma; \mathbf{z} - \boldsymbol{\eta}, \boldsymbol{\eta}) &\leq u'(\mathbf{y}^*, \gamma; -\boldsymbol{\xi}, \boldsymbol{\eta}) - u'(\mathbf{y}^*, \gamma; -\boldsymbol{\xi}, \boldsymbol{\xi}) + u'(\mathbf{y}^*, \gamma; -\boldsymbol{\xi}, \boldsymbol{\xi}) - C(\boldsymbol{\xi} - \boldsymbol{\eta}) + C(-\mathbf{z}) \\ &\leq u'(\mathbf{y}^*, \gamma; \mathbf{0}, \boldsymbol{\eta} - \boldsymbol{\xi}) + u'(\mathbf{y}^*, \gamma; -\boldsymbol{\xi}, \boldsymbol{\xi}) - C(\boldsymbol{\xi} - \boldsymbol{\eta}) + C(-\mathbf{z}) \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{\rho}{2} c_{\max} \sum_{i=1}^n |\xi_i - \eta_i| + \sum_{i=1}^n \beta_i \xi_i - \frac{c_{\min}}{2} \sum_{i=1}^n |\xi_i - \eta_i| + C(-\mathbf{z}) \\
 &\leq \sum_{i=1}^n \beta_i \eta_i + \frac{\rho c_{\max} - c_{\min}}{2} \sum_{i=1}^n |\xi_i - \eta_i| + C(-\mathbf{z}) \\
 &\leq \sum_{i=1}^n (\beta_i + \rho c_{\max} - c_{\min}) \eta_i + C(-\mathbf{z}),
 \end{aligned}$$

where the second inequality follows by Lemma 5 and the third inequality follows by Lemma 4.

Now we show equation (34). Suppose $(\mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta})$ is a feasible direction. Again, let \mathbf{y}^* be an optimal solution to equation (16) at (\mathbf{x}, γ) . Then $\mathbf{y}^* + t\boldsymbol{\eta}$ is clearly a feasible solution to equation (16) at $(\mathbf{x} + t\mathbf{z} + t\boldsymbol{\eta}, \gamma - t\boldsymbol{\eta})$ and thus we have

$$\begin{aligned}
 v'(\mathbf{x}, \gamma; \mathbf{z} + \boldsymbol{\eta}, -\boldsymbol{\eta}) &= \lim_{t \rightarrow 0} \frac{v(\mathbf{x} + t\mathbf{z} + t\boldsymbol{\eta}, \gamma - t\boldsymbol{\eta}) - v(\mathbf{x}, \gamma)}{t} \\
 &\leq \lim_{t \rightarrow 0} \frac{u(\mathbf{y}^* + t\boldsymbol{\eta}, \gamma - t\boldsymbol{\eta}) + C(\mathbf{y}^* - \mathbf{x} - t\mathbf{z}) - u(\mathbf{y}^*, \gamma) - C(\mathbf{y}^* - \mathbf{x})}{t} \\
 &\leq u'(\mathbf{y}^*, \gamma; \boldsymbol{\eta}, -\boldsymbol{\eta}) + C(-\mathbf{z}) \\
 &\leq \sum_{i=1}^n \beta_i \eta_i + C(-\mathbf{z}),
 \end{aligned}$$

where the second inequality follows from the subadditivity and the positive homogeneity of $C(\cdot)$ and the last inequality follows from the assumption.

To show the result in (33), suppose $(\mathbf{0}, \mathbf{z})$ is a feasible direction at (\mathbf{x}, γ) . Let \mathbf{y}^* be an optimal solution to equation (16) at (\mathbf{x}, γ) . Then \mathbf{y}^* is a feasible solution for $(\mathbf{x}, \gamma + t\mathbf{z})$ and thus,

$$\begin{aligned}
 v'(\mathbf{x}, \gamma; \mathbf{0}, \mathbf{z}) &= \lim_{t \rightarrow 0} \frac{v(\mathbf{x}, \gamma + t\mathbf{z}) - v(\mathbf{x}, \gamma)}{t} \\
 &\leq \lim_{t \rightarrow 0} \frac{u(\mathbf{y}^*, \gamma + t\mathbf{z}) + C(\mathbf{x} - \mathbf{y}^*) - u(\mathbf{y}^*, \gamma) - C(\mathbf{x} - \mathbf{y}^*)}{t} \\
 &= \lim_{t \rightarrow 0} \frac{u(\mathbf{y}^*, \gamma + t\mathbf{z}) - u(\mathbf{y}^*, \gamma)}{t} = u'(\mathbf{y}^*, \gamma; \mathbf{0}, \mathbf{z}) \leq \frac{\rho c_{\max}}{2} \sum_{i=1}^n |z_i|,
 \end{aligned}$$

which completes the proof.

The last piece of Theorem 1 is the Lipschitz continuity of $u_t(\cdot)$. In the next proposition, we show that item (1) and (1) imply item (1).

LEMMA 8. Consider a bounded function $f: \Delta \rightarrow \mathbb{R}$ that is convex, continuous, and satisfies

1. $|f'(\mathbf{y}, \gamma; \mp \boldsymbol{\eta}, \pm \boldsymbol{\eta})| \leq \sum_{i=1}^n \beta_i \eta_i$ for all $(\mathbf{x}, \gamma) \in \Delta$ and any feasible direction $(\mp \boldsymbol{\eta}, \pm \boldsymbol{\eta})$ with $\boldsymbol{\eta} \geq \mathbf{0}$;
2. $f'(\mathbf{y}, \gamma; \mathbf{0}, \mathbf{v}) \leq (\rho c_{\max}/2) \sum_{i=1}^n |v_i|$ for all $(\mathbf{x}, \gamma) \in \Delta$ and any feasible direction $(\mathbf{0}, \mathbf{v})$ with $\mathbf{e}^T \mathbf{v} = 0$.

Then, the function $f(\cdot)$ is Lipschitz continuous on Δ° with Lipschitz constant $(3/2) \sqrt{2n} \beta_i$.

Proof. Let $(\mathbf{y}, \gamma) \in \Delta^\circ$. Let $\boldsymbol{\xi}^+ = \max(\boldsymbol{\xi}, 0)$ and $\boldsymbol{\xi}^- = \min(\boldsymbol{\xi}, 0)$, then $(\boldsymbol{\xi}, \boldsymbol{\eta}) = (\boldsymbol{\xi}^+, -\boldsymbol{\xi}^+) + (\boldsymbol{\xi}^-, -\boldsymbol{\xi}^-) + (0, \boldsymbol{\xi}^+ + \boldsymbol{\xi}^- + \boldsymbol{\eta})$ and it follows from Theorem 2 and Assumption 2 that

$$\begin{aligned}
 |f'(\mathbf{y}, \gamma; \boldsymbol{\xi}, \boldsymbol{\eta})| &\leq |f'(\mathbf{y}, \gamma; \boldsymbol{\xi}^+, -\boldsymbol{\xi}^+)| + |f'(\mathbf{y}, \gamma; \boldsymbol{\xi}^-, -\boldsymbol{\xi}^-)| + |f'(\mathbf{y}, \gamma; 0, \boldsymbol{\xi}^+ + \boldsymbol{\xi}^- + \boldsymbol{\eta})| \\
 &\leq \sum_{i=1}^n \beta_i |\xi_i^+| + \sum_{i=1}^n \beta_i |\xi_i^-| + (\rho c_{\max}/2) \|\boldsymbol{\xi}^+ + \boldsymbol{\xi}^- + \boldsymbol{\eta}\|_1 \\
 &\leq \beta_{\max} \|\boldsymbol{\xi}\|_1 + (\rho c_{\max}/2) \|\boldsymbol{\xi}^+ + \boldsymbol{\xi}^- + \boldsymbol{\eta}\|_1 \\
 &\leq \beta_{\max} \|\boldsymbol{\xi}\|_1 + (\beta_{\max}/2) (\|\boldsymbol{\xi}^+\|_1 + \|\boldsymbol{\xi}^-\|_1 + \|\boldsymbol{\eta}\|_1) \\
 &\leq (3/2) \beta_{\max} \|(\boldsymbol{\xi}, \boldsymbol{\eta})\|_1 \leq (3/2) \sqrt{2n} \beta_{\max} \|(\boldsymbol{\xi}, \boldsymbol{\eta})\|_2.
 \end{aligned}$$

With all the pieces ready, the proofs for Theorem 1 for $(u_t(\cdot))$ and $(v_t(\cdot))$ thus follow from Proposition 3, Proposition 4, Proposition 5, and the induction.

A.4. Other Proofs

Proof of Theorem 3: Fix $\gamma \in S$. Let $\mathbf{y}^*(\mathbf{x}, \gamma) = \{\mathbf{y} \in \Delta_{n-1}(I) : u(\mathbf{x}, \gamma) = C(\mathbf{y} - \mathbf{x}) + u(\mathbf{y}, \gamma)\}$ be the set of optimal solutions corresponding to the system state $\mathbf{x} \in S$. It is easy to verify that

$$\Omega_u(\gamma) = \cup_{\mathbf{x} \in \Delta_{n-1}(I)} \mathbf{y}^*(\mathbf{x}, \gamma). \quad (36)$$

As $C(\cdot)$ and $u(\cdot)$ are continuous and $\Delta_{n-1}(I)$ is compact, by Berge's Maximum Theorem, $\mathbf{y}^*(\cdot)$ is a nonempty-valued and compact-valued upper hemicontinuous¹⁶ correspondence. As $C(\cdot)$ and $u(\cdot)$ are also convex, $\mathbf{y}^*(\cdot)$ is also convex-valued. So, it is clear from (36) that $\Omega_u(\gamma)$ is nonempty. To show $\Omega_u(\gamma)$ is compact, suppose $\mathbf{y}_1, \mathbf{y}_2, \dots$ is a sequence in $\Omega_u(\gamma)$ such that $\mathbf{y}_n \in \mathbf{y}^*(\mathbf{x}_n, \gamma)$ for $n \in \mathbb{N}$ and $\mathbf{y}_n \rightarrow \mathbf{y}$. We need to show that $\mathbf{y} \in \Omega_u(\gamma)$. By passing through a subsequence, we may assume that $\mathbf{y}_{n_k} \in \mathbf{y}^*(\mathbf{x}_{n_k}, \gamma)$, $\mathbf{x}_{n_k} \rightarrow \mathbf{x}$ and $\mathbf{y}_{n_k} \rightarrow \mathbf{y}$. As $\mathbf{y}^*(\cdot)$ is compact-valued, by the Closed Graph Theorem, $\mathbf{y}^*(\cdot)$ has a closed graph. This implies that $\mathbf{y} \in \mathbf{y}^*(\mathbf{x}, \gamma) \subset \Omega_u(\gamma)$, and therefore $\Omega_u(\gamma)$ is compact.

To show that $\Omega_u(\gamma)$ is connected, suppose the reverse is true. Then, there exist open sets V_1, V_2 in $\Delta_{n-1}(I)$ such that $V_1 \cap V_2 = \emptyset$, $V_1 \cup V_2 \supset \Omega_u(\gamma)$, and $V_1 \cap \Omega_u(\gamma)$ and $V_2 \cap \Omega_u(\gamma)$ are nonempty. As $\mathbf{y}^*(\cdot)$ is convex-valued, this implies that, for any $\mathbf{x} \in \Delta_{n-1}(I)$, $\mathbf{y}^*(\mathbf{x}, \gamma)$ is either in V_1 or in V_2 , but not both. Let $(U_1, \gamma) = \mathbf{y}^{*-1}(V_1)$ and $(U_2, \gamma) = \mathbf{y}^{*-1}(V_2)$. Then U_1, U_2 are open, $U_1 \cap U_2 = \emptyset$, $U_1 \cup U_2 \supset \Delta_{n-1}(I)$, and $U_1 \cap \Delta_{n-1}(I)$ and $U_2 \cap \Delta_{n-1}(I)$ are nonempty. This implies that the $(n-1)$ -dimensional simplex $\Delta_{n-1}(I)$ is not connected. We have reached a contradiction. Therefore, $\Omega_u(\gamma)$ is also connected.

Next, to show that π^* is optimal, note that $\pi^*(\mathbf{x}, \gamma) = \mathbf{x}$ for $\mathbf{x} \in \Omega_u(\gamma)$ is clear from (17). If $\mathbf{x} \notin \Omega_u(\gamma)$, then, by (36), $\pi^*(\mathbf{x}, \gamma) \in \Omega_u(\gamma)$. Now, suppose there exists $\pi^*(\mathbf{x}, \gamma) = \mathbf{y} \in \Omega_u(\gamma)^\circ$, then $\mathbf{y} + t(\mathbf{x} - \mathbf{y}) \in \Omega_u(\gamma)$ for small enough $t > 0$. Set $\mathbf{z} = \mathbf{y} + t(\mathbf{x} - \mathbf{y})$. Then $u(\mathbf{z}, \gamma) + C(\mathbf{z} - \mathbf{x}) \leq u(\mathbf{y}, \gamma) + C(\mathbf{y} - \mathbf{z}) + C(\mathbf{z} - \mathbf{x}) = u(\mathbf{y}, \gamma) + tC(\mathbf{y} - \mathbf{x}) + (1-t)C(\mathbf{y} - \mathbf{x}) = u(\mathbf{y}, \gamma) + C(\mathbf{y} - \mathbf{x})$. So, \mathbf{z} is as good a solution as \mathbf{y} . Therefore, there exists an optimal solution $\pi^*(\mathbf{x}, \gamma) \in \mathcal{B}(\Omega_u(\gamma))$ if $\mathbf{x} \notin \Omega_u(\gamma)$. \square

Proof of Proposition 1: Suppose $\mathbf{x} \in \Omega_u(\gamma)$. Take any feasible direction $(\mathbf{z}, \mathbf{0})$ at (\mathbf{x}, γ) . Then, by (17),

$$\frac{u(\mathbf{x} + t\mathbf{z}, \gamma) - u(\mathbf{x}, \gamma)}{t} \geq -C(\mathbf{z})$$

for $t > 0$. Taking the limit as $t \downarrow 0$, we have $u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0}) \geq -C(\mathbf{z})$. Conversely, suppose $u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0}) \geq -C(\mathbf{z})$ for any feasible direction \mathbf{z} at \mathbf{x} in \mathcal{H} . Let $\phi(t) = u(\mathbf{x} + t\mathbf{z}, \gamma)$. Then, $\phi(\cdot)$ is convex, $\phi(0) = u(\mathbf{x})$, and $\phi'(0+) = u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0}) \geq -C(\mathbf{z})$. By the subgradient inequality, $t\phi'(0+) + \phi(0) \leq \phi(t)$. This implies that $-tC(\mathbf{z}) + u(\mathbf{x}, \gamma) \leq u(\mathbf{x} + t\mathbf{z}, \gamma)$ is true for any feasible direction $(\mathbf{z}, \mathbf{0})$. Therefore, we have $\mathbf{x} \in \Omega_u(\gamma)$. \square

¹⁶ Upper hemicontinuity can be defined as follows. Suppose X and Y are topological spaces. A correspondence $f: X \rightarrow \mathcal{P}(Y)$ (power set of Y) is upper hemicontinuous if for any open set V in Y , $f^{-1}(V) = \{\mathbf{x} \in X | f(\mathbf{x}) \subset V\}$ is open in X .

Proof of Proposition 2: For the “if” part, suppose $\mathbf{x} \notin \Omega_u(\gamma)$. Then, there exists $\mathbf{y} \in \Delta_{n-1}(I)$ such that $u(\mathbf{x}, \gamma) > C(\mathbf{y} - \mathbf{x}) + u(\mathbf{y}, \gamma)$. Take any $\mathbf{g} \in \partial_x u(\mathbf{x}, \gamma)$. By the subgradient inequality, $u(\mathbf{x}, \gamma) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}) \leq u(\mathbf{y}, \gamma)$. It follows that

$$C(\mathbf{y} - \mathbf{x}) < -\mathbf{g}^T(\mathbf{y} - \mathbf{x}).$$

Suppose $\mathbf{w} = (w_{ij})$ is an optimal solution to problem (1). Then $C(\mathbf{y} - \mathbf{x}) = \sum_i \sum_j c_{ij} w_{ij}$, and by Lemma 3, $-\mathbf{g}^T(\mathbf{y} - \mathbf{x}) = \sum_i g_i (y_i - x_i)^- - \sum_j g_j (y_j - x_j)^+ = \sum_i \sum_j (g_i - g_j) w_{ij}$. So, we have

$$\sum_i \sum_j c_{ij} z_{ij} < \sum_i \sum_j (g_i - g_j) w_{ij}.$$

Hence, there exists i and j such that $g_i - g_j > c_{ij}$. This implies $\mathbf{g} \notin \mathcal{G}$.

For the “only if” part, suppose $\mathbf{x} > \mathbf{0}$ and $\mathbf{x} \in \Omega_u(\gamma)$. Assume $\partial_x u(\mathbf{x}, \gamma) \cap \mathcal{G} = \emptyset$. We will show that this leads to a contradiction. Let P be the orthogonal projection from \mathbb{R}^n to the subspace $\mathcal{H} = \{\mathbf{x} \in \mathbb{R}^n : \sum_i x_i = 0\}$. Then

$$P(\mathbf{x}) = \mathbf{x} - \frac{\sum_i x_i}{n} \mathbf{e},$$

where $\mathbf{e} = (1, \dots, 1)$ in \mathbb{R}^n . Noting that $\mathcal{G} + \alpha \mathbf{e} \subset \mathcal{G}$ for any $\alpha \in \mathbb{R}$, it is easy to verify that

$$\partial_x u(\mathbf{x}, \gamma) \cap \mathcal{G} = \emptyset \quad \text{if and only if} \quad \partial_x u(\mathbf{x}, \gamma) \cap P(\mathcal{G}) = \emptyset,$$

since $\partial_x u(\mathbf{x}, \gamma) \subseteq \mathcal{H}$. As $\partial_x u(\mathbf{x}, \gamma)$ is closed and $P(\mathcal{G})$ is compact, by Hahn-Banach Theorem, there exists $\mathbf{z} \in \mathcal{H}$, $a \in \mathbb{R}$ and $b \in \mathbb{R}$ such that

$$\langle \mathbf{g}, \mathbf{z} \rangle < a < b < \langle \boldsymbol{\lambda}, \mathbf{z} \rangle$$

for every $\mathbf{g} \in P(\partial_x u(\mathbf{x}, \gamma))$ and for every $\boldsymbol{\lambda} \in P(\mathcal{G})$, or equivalently, as $\langle \mathbf{g}, \mathbf{z} \rangle = \langle P(\mathbf{g}), \mathbf{z} \rangle$ and $\langle \boldsymbol{\lambda}, \mathbf{z} \rangle = \langle P(\boldsymbol{\lambda}), \mathbf{z} \rangle$, for every $\mathbf{g} \in \partial_x u(\mathbf{x}, \gamma)$ and for every $\boldsymbol{\lambda} \in \mathcal{G}$. As \mathbf{z} is a feasible direction in \mathcal{H} at $\mathbf{x} \in \Omega_u(\gamma)$, by Proposition 1, we have $u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0}) \geq -C(\mathbf{z})$. It follows that

$$\sup\{\langle \mathbf{g}, \mathbf{z} \rangle : \mathbf{g} \in \partial_x u(\mathbf{x}, \gamma)\} = u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0}) \geq -C(\mathbf{z}).$$

So, we have

$$-C(\mathbf{z}) \leq a < b < \langle \boldsymbol{\lambda}, \mathbf{z} \rangle$$

for every $\boldsymbol{\lambda} \in \mathcal{G}$. However, by the dual formulation (23), there exists $\boldsymbol{\lambda} \in \{(y_1, y_2, \dots, y_n) \mid y_j - y_i \leq c_{ij} \forall i, j\}$ such that $\langle \boldsymbol{\lambda}, \mathbf{z} \rangle = C(\mathbf{z})$, or equivalently, $\langle -\boldsymbol{\lambda}, \mathbf{z} \rangle = -C(\mathbf{z})$. Recognizing that $-\boldsymbol{\lambda} \in \mathcal{G}$ leads to the contradiction. Therefore, it follows that $\partial_x u(\mathbf{x}, \gamma) \cap \mathcal{G} \neq \emptyset$. \square

Proof of Corollary 2: If $u(\cdot, \gamma)$ is differentiable at \mathbf{x} , then

$$\partial_x u(\mathbf{x}, \gamma) = \left\{ \left(\frac{\partial u(\mathbf{x}, \gamma)}{\partial x_1}, \frac{\partial u(\mathbf{x}, \gamma)}{\partial x_2}, \dots, \frac{\partial u(\mathbf{x}, \gamma)}{\partial x_n} \right) \right\}.$$

In this case, it is easy to see that (??) is simplified to (??). To show that $\mathbf{x} \in \Omega_u(\gamma)$ implies (??) for $\mathbf{x} \in \mathcal{B}(\Delta_{n-1}(I))$. Note that the equality $\sup\{\mathbf{g}^T \mathbf{z} : \mathbf{g} \in \partial_x u(\mathbf{x}, \gamma)\} = u'(\mathbf{x}, \gamma; \mathbf{z}, \mathbf{0})$ now holds for $\mathbf{x} \in \mathcal{B}(\Delta_{n-1}(I))$. The rest of the proof is the same as Proposition 2. \square

Proof of Theorem 2: To show that value function retains its structure in the infinite horizon setting, we invoke the general approach outlined in Porteus (1975) and Porteus (1982) which “iterates” the structural properties of the one-stage problem. Let \mathcal{V}^* be the space of convex, continuous and bounded functions over Δ . Note that a one-step *structure preservation* property holds by Lemma 3, Lemma 4, and Lemma 5: combined, they say that if the next period value function is in \mathcal{V}^* , then the optimal value of the current period is also in \mathcal{V}^* . Furthermore, the set \mathcal{V}^* with the sup-norm $\|\cdot\|_\infty$ is a complete metric space. These two observations allow us to apply Corollary 1 of Porteus (1975) and conclude that $v \in \mathcal{V}^*$ (the remaining assumptions needed to apply the result can be easily checked). The rest of the proof follows from Lemma 3, Lemma 4, Lemma 5, Theorem 3, Proposition 1, Proposition 2, and Corollary 2. \square

Proof of Lemma 1: If $\mathbf{x} \in D^k$, $\mathbf{a}_k \in \partial_x u_J(\mathbf{x}, \boldsymbol{\gamma})$. Since $a_{ki} - a_{kj} \leq c_{ij}$ for all i, j , according to Proposition 2, we have $\mathbf{x} \in \Omega_{u_J}(\boldsymbol{\zeta})$. For the second part, we first write down the primal formulation for problem (??):

$$\begin{aligned} \bar{v}_J(\mathbf{x}, \boldsymbol{\zeta}) = \quad & \min \quad \mathbf{c} \cdot \mathbf{w} + \xi \\ & \text{subject to} \quad \sum_{i=1}^n w_{i,j} - \sum_{k=1}^n w_{j,k} = z_j - x_j \quad \forall j = 1, 2, \dots, n \\ & \mathbf{w} \geq 0 \\ & \mathbf{e}^T \mathbf{z} = \mathbf{e}^T \mathbf{x} \\ & \xi \geq (\mathbf{z} - \mathbf{y}_k)^T \mathbf{a}_k + (\boldsymbol{\zeta} - \boldsymbol{\gamma}_k)^T \mathbf{b}_k + \boldsymbol{\nu}_k \quad \forall k = 1, 2, \dots, J \\ & \mathbf{z} \geq 0. \end{aligned}$$

Since $\mathbf{x} \in \Omega_{u_J}(\boldsymbol{\zeta})$, one optimal solution to the primal formulation is $\mathbf{w} = 0$, $\mathbf{z} = \mathbf{x}$, $\xi = (\mathbf{x} - \mathbf{y}_k)^T \mathbf{a}_k + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_k)^T \mathbf{b}_k + \boldsymbol{\nu}_k$. The dual solution $\boldsymbol{\lambda} = \mathbf{a}_k$, $\mu_k = 1$, $\mu_l = 0$, $\forall l \neq k$ is clearly feasible. It also satisfies the complementary slackness condition. Therefore, the solution is optimal. \square

Proof of Theorem 4: Let us first introduce some notation. For any bounded function $f : \Delta \rightarrow \mathbb{R}$, we define the mapping \mathcal{L} such that $\mathcal{L}f : \Delta \rightarrow \mathbb{R}$ is another bounded function such that

$$(\mathcal{L}f)(\mathbf{y}, \boldsymbol{\gamma}) = l(\mathbf{y}) + \rho \int \min_{\mathbf{y}' \in \Delta_{n-1}(\mathbf{e}^T \mathbf{x}')} C(\mathbf{y}' - \mathbf{x}') + f(\mathbf{y}', \boldsymbol{\gamma}') d\boldsymbol{\mu} \quad \forall (\mathbf{y}, \boldsymbol{\gamma}) \in \Delta, \quad (37)$$

where $\mathbf{x}' = \tau_x(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}, \mathbf{P})$ and $\boldsymbol{\gamma}' = \tau_\gamma(\mathbf{y}, \boldsymbol{\gamma}, \mathbf{d}, \mathbf{P})$. Note that \mathcal{L} is closely related to the standard Bellman operator associated with the MDP defined in (9); see, for example, Bertsekas and Tsitsiklis (1996). The difference from the standard definition is that \mathcal{L} comes from the Bellman recursion for $u(\mathbf{y}, \boldsymbol{\gamma})$ instead of $v(\mathbf{x}, \boldsymbol{\gamma})$. With this in mind, we henceforth simply refer to \mathcal{L} as the “Bellman operator” and note a few standard properties.

LEMMA 9. *The Bellman operator \mathcal{L} has the following properties:*

1. (Monotonicity) *Given bounded $f_1, f_2 : \Delta \rightarrow \mathbb{R}$ with $f_1 \leq f_2$, then $\mathcal{L}f_1 \leq \mathcal{L}f_2$.*
2. (Contraction) *For any bounded $f_1, f_2 : \Delta \rightarrow \mathbb{R}$, it holds that $\|\mathcal{L}f_1 - \mathcal{L}f_2\|_\infty \leq \rho \|f_1 - f_2\|_\infty$.*
3. (Fixed Point) *The optimal value function u is the unique fixed point of \mathcal{L} , i.e., $\mathcal{L}u = u$.*
4. (Constant Shift) *Let $\mathbf{1}$ be the constant one function, i.e., $\mathbf{1}(\cdot) = 1$, and let α be a scalar. For any bounded $f : \Delta \rightarrow \mathbb{R}$, it holds that $\mathcal{L}(f + \alpha\mathbf{1}) = \mathcal{L}f + \rho\alpha\mathbf{1}$.*

Proof. These basic properties for our Bellman operator \mathcal{L} are well-known for the standard Bellman operator and can be proved in an analogous manner; see, for example, Puterman (1994) or Bertsekas and Tsitsiklis (1996). \square

We now move on to the main proof. We want to show that for each $\epsilon > 0$, there exists an almost surely finite iteration index $J(\epsilon)$ such that for all $J \geq J(\epsilon)$, it holds that $\|u_J - u\|_\infty \leq \epsilon$. Let $B_r(\mathbf{y}, \boldsymbol{\gamma})$ be a $(2n - 1)$ -dimensional ball centered at $(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta^\circ$ with radius r . Consider some $\epsilon' > 0$ (to be specified later) and let $\mathcal{C}(\epsilon')$ be an ϵ' -covering of Δ , meaning that $\mathcal{C}(\epsilon')$ is a finite collection of points in Δ° (representing the centers of a finite collection of balls with radius ϵ') and $\Delta \subseteq \bigcup_{(\mathbf{y}, \boldsymbol{\gamma}) \in \mathcal{C}(\epsilon')} B_{\epsilon'}(\mathbf{y}, \boldsymbol{\gamma})$. Let $(\tilde{\mathbf{y}}_1, \tilde{\boldsymbol{\gamma}}_1), (\tilde{\mathbf{y}}_2, \tilde{\boldsymbol{\gamma}}_2), \dots$ denote the sequence of sample points visited by the algorithm (one per iteration). Thus, by Assumption 3, we have $\sum_J \mathbf{P}\{(\tilde{\mathbf{y}}_J, \tilde{\boldsymbol{\gamma}}_J) \in B_{\epsilon'}(\mathbf{y}, \boldsymbol{\gamma})\} = \infty$, and an application of the Borel-Cantelli lemma tells us that each ball $B_{\epsilon'}(\mathbf{y}, \boldsymbol{\gamma})$ associated with the covering is visited infinitely often with probability one. To reduce notation, we will often suppress $(\mathbf{y}, \boldsymbol{\gamma})$ and use $B_{\epsilon'}$ to denote a generic ball in the covering. Our proof follows three main ideas:

1. For any infinite trajectory of sampled states, we can split it into an infinite number of “phases” such that in each phase, every ball associated with the ϵ' -covering is visited at least once.

2. We can then construct an auxiliary “batch” algorithm whose iteration counter aligns with the sequence of phases from the previous step. This new algorithm is defined as another instance of Algorithm 1, where on any given iteration, we group all states visited in the corresponding phase of the main algorithm into a single batch and perform all updates at once. For clarity, we will refer to the main algorithm as the “asynchronous” version of the batch algorithm.

3. The auxiliary batch algorithm can be viewed as an approximate version of value iteration. Using the properties of \mathcal{L} , we can show that it converges to an approximation of u (with error depending on ϵ'). Finally, we conclude by arguing that the main algorithm does not deviate too far from the auxiliary version.

Let $J_0 = 0$ and for $K = 1, 2, \dots$, define the random variable

$$J_{K+1} = \min\{J > J_K : \forall (\mathbf{y}, \boldsymbol{\gamma}) \in \mathcal{C}(\epsilon'), \exists J' \text{ s.t. } J_K < J' \leq J, (\tilde{\mathbf{y}}_{J'}, \tilde{\boldsymbol{\gamma}}_{J'}) \in B_{\epsilon'}(\mathbf{y}, \boldsymbol{\gamma})\}$$

to be the first time after J_K such that every ball in the ϵ' -covering is visited at least once. Notably, J_1 is the first time that the entire covering is visited at least once. We denote the set of iterations

$$\mathcal{J}_K = \{J_{K-1} + 1, J_{K-1} + 2, \dots, J_K\}$$

to be the “ K th phase” of the algorithm and let $\hat{\mathcal{S}}_K = \{(\tilde{\mathbf{y}}_J, \tilde{\boldsymbol{\gamma}}_J)\}_{J \in \mathcal{J}_K}$ be the set of states visited throughout the course of phase K .

We now describe “path-dependent” instances of Algorithm 1 to assist with the remaining analysis. To be precise with the definitions, let us consider a sample path ω . The auxiliary batch algorithm associated with ω is a new instance of Algorithm 1 that uses iteration counter K and generates hyperplanes at the set of states $\mathcal{S}_K = \hat{\mathcal{S}}_K(\omega)$ for all $K \geq 1$. The initial approximation is $\hat{u}_0 = u_0$ and the estimate after K batch updates is denoted $\hat{u}_K(\mathbf{y}, \boldsymbol{\gamma})(\omega) = \max_{i=1, \dots, N_K} \hat{g}_i(\mathbf{y}, \boldsymbol{\gamma})(\omega)$. We are now interested in studying the stochastic process $\{\hat{u}_K(\mathbf{y}, \boldsymbol{\gamma})\}$.

Next, we observe that the hyperplanes generated at iteration $K + 1$ of the batch algorithm are tangent to $\mathcal{L}\hat{u}_K$ at the points in \mathcal{S}_{K+1} . Let $\kappa = (3/2)\sqrt{2n}\beta$. Note that by repeatedly applying Lemma 8, Proposition 5, Proposition 3, and Proposition 4, and using $u_0 = 0$, we can argue that all (tangent) hyperplanes generated throughout the algorithm have directional derivatives bounded by κ . It follows that if $(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\gamma}})$ is a sample point in \mathcal{S}_{K+1} that lies in a ball $B_{\epsilon'}$ and it generates a hyperplane \hat{g} , then the underestimation error within the ball is upper-bounded by $\max_{(\mathbf{y}, \boldsymbol{\gamma}) \in B_{\epsilon'}} [(\mathcal{L}\hat{u}_K)(\mathbf{y}, \boldsymbol{\gamma}) - \hat{g}(\mathbf{y}, \boldsymbol{\gamma})] \leq 2\kappa\epsilon'$ (using the fact that there is zero estimation error at $(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\gamma}})$, the tangent point). Applying this across the ϵ' -covering, we have:

$$\mathcal{L}\hat{u}_K - (2\kappa\epsilon')\mathbf{1} \leq \max_{i=N_K+1, \dots, N_{K+1}} \hat{g}_i \leq \max_{i=1, \dots, N_{K+1}} \hat{g}_i = \hat{u}_{K+1}. \quad (38)$$

Therefore, we have a form of approximate value iteration and can analyze it accordingly (see Bertsekas and Tsitsiklis (1996)). Utilizing the monotonicity and shift properties of Lemma 9, we apply \mathcal{L} to both sides of (38) for $K = 0$ to obtain

$$\mathcal{L}(\mathcal{L}\hat{u}_0 - 2\kappa\epsilon'\mathbf{1}) = \mathcal{L}^2\hat{u}_0 - \rho(2\kappa\epsilon')\mathbf{1} \leq \mathcal{L}\hat{u}_1.$$

Subtracting $(2\kappa\epsilon')\mathbf{1}$ from both sides and then applying (38) for $K = 1$, we have

$$\mathcal{L}^2\hat{u}_0 - \rho(2\kappa\epsilon')\mathbf{1} - (2\kappa\epsilon')\mathbf{1} \leq \mathcal{L}\hat{u}_1 - (2\kappa\epsilon')\mathbf{1} \leq \hat{u}_2.$$

Iterating these steps, we see that $\mathcal{L}^K\hat{u}_0 - (2\kappa\epsilon')(1 + \rho + \dots + \rho^{K-1})\mathbf{1} \leq \hat{u}_K$. Taking limits, using the convergence of the value iteration algorithm (see Puterman (1994)), and noting that $\hat{u}_K \leq u$ for all K , we arrive at

$$u(\mathbf{y}, \boldsymbol{\gamma}) - \frac{2\kappa\epsilon'}{1-\rho} \leq \lim_{K \rightarrow \infty} \hat{u}_K(\mathbf{y}, \boldsymbol{\gamma}) \leq u(\mathbf{y}, \boldsymbol{\gamma}), \quad \forall (\mathbf{y}, \boldsymbol{\gamma}) \in \Delta. \quad (39)$$

Hence, we have shown that the auxiliary batch algorithm generates value function approximations that closely approximate u in the limit.

The final step is to relate the main asynchronous algorithm to the auxiliary batch version. We claim that the value function approximation \hat{u}_K generated by the K th phase, for $K \geq 1$, of the batch algorithm is within a certain error bound of the approximation from the asynchronous algorithm at J_K :

$$\hat{u}_K - (4\kappa\epsilon')(1 + \rho + \dots + \rho^{K-1})\mathbf{1} \leq u_{J_K}. \quad (40)$$

To prove (40), let us consider the first phase, $K = 1$. Recall that the two algorithms are initialized with identical approximations, so $\hat{u}_0 = u_0$. Since $\{u_J\}$ is a nondecreasing sequence of functions, we have $\hat{u}_0 \leq u_J$ and $\mathcal{L}\hat{u}_0 \leq \mathcal{L}u_J$ for any $J \in \mathcal{J}_1$ by the monotonicity property of Lemma 9. Also note that the auxiliary batch algorithm builds a uniform underestimate of $\mathcal{L}\hat{u}_0$ with points of tangency belonging to $\hat{\mathcal{S}}_1$, so we have $\hat{u}_1 \leq \mathcal{L}\hat{u}_0 \leq \mathcal{L}u_J$. The hyperplane g_{J+1} added in iteration $J + 1 \in \mathcal{J}_1$ of the asynchronous algorithm is tangent to $\mathcal{L}u_J$ at $(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1})$, so it follows that

$$\hat{u}_1(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1}) \leq (\mathcal{L}u_J)(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1}) = g_{J+1}(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1}).$$

Suppose $(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1})$ is in a ball $B_{\epsilon'}$. Then g_{J+1} can fall below \hat{u}_1 by at most $\max_{(\mathbf{y}, \boldsymbol{\gamma}) \in B_{\epsilon'}} [\hat{u}_1(\mathbf{y}, \boldsymbol{\gamma}) - g_{J+1}(\mathbf{y}, \boldsymbol{\gamma})] \leq 4\kappa\epsilon'$ within the ball. Since this holds for every hyperplane (and corresponding ball) added

throughout the phase $K = 1$ and noting that every point in Δ can be associated with some well-approximating hyperplane (due to the property that each phase contains at least one visit to every ball), we have

$$\hat{u}_1 - (4\kappa\epsilon') \mathbf{1} \leq \max_{J=1, \dots, J_1} g_J = u_{J_1}, \quad (41)$$

which proves (40) for $K = 1$. Applying \mathcal{L} to both sides of (41), utilizing Lemma 9, noting that \hat{u}_2 underestimates $\mathcal{L}\hat{u}_1$, and applying the nondecreasing property of the sequence $\{u_J\}$, we have

$$\hat{u}_2 - \rho(4\kappa\epsilon') \mathbf{1} \leq \mathcal{L}\hat{u}_1 - \rho(4\kappa\epsilon') \mathbf{1} \leq \mathcal{L}u_{J_1} \leq \mathcal{L}u_J, \quad \forall J \in \mathcal{J}_2.$$

By the same reasoning as above, for $J + 1 \in \mathcal{J}_2$, it must hold that

$$\hat{u}_2(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1}) - \rho(4\kappa\epsilon') \leq (\mathcal{L}u_J)(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1}) = g_{J+1}(\tilde{\mathbf{y}}_{J+1}, \tilde{\boldsymbol{\gamma}}_{J+1}),$$

and we obtain $\hat{u}_2 - \rho(4\kappa\epsilon') \mathbf{1} - (4\kappa\epsilon') \mathbf{1} \leq u_{J_2}$, proving (40) for $K = 2$. We can iterate these steps to argue (40) for any K . Taking limits (all subsequent limits exist due to the boundedness and monotonicity of the sequences), we get

$$\lim_{K \rightarrow \infty} \hat{u}_K(\mathbf{y}, \boldsymbol{\gamma}) - \frac{4\kappa\epsilon'}{1-\rho} \leq \lim_{K \rightarrow \infty} u_{J_K}(\mathbf{y}, \boldsymbol{\gamma}) = \lim_{J \rightarrow \infty} u_J(\mathbf{y}, \boldsymbol{\gamma}) \leq u(\mathbf{y}, \boldsymbol{\gamma}), \quad \forall (\mathbf{y}, \boldsymbol{\gamma}) \in \Delta, \quad (42)$$

where the equality follows by the fact $u_{J_K}(\mathbf{y}, \boldsymbol{\gamma})$ is a subsequence of $u_J(\mathbf{y}, \boldsymbol{\gamma})$. Combining (39) and (42), we obtain

$$u(\mathbf{y}, \boldsymbol{\gamma}) - \frac{6\kappa\epsilon'}{1-\rho} \leq \lim_{K \rightarrow \infty} \hat{u}_K(\mathbf{y}, \boldsymbol{\gamma}) - \frac{4\kappa\epsilon'}{1-\rho} \leq \lim_{J \rightarrow \infty} u_J(\mathbf{y}, \boldsymbol{\gamma}) \leq u(\mathbf{y}, \boldsymbol{\gamma}), \quad \forall (\mathbf{y}, \boldsymbol{\gamma}) \in \Delta,$$

showing that the asynchronous algorithm generates value function approximations that are arbitrarily close to u : if we set $\epsilon' = (1-\rho)\epsilon/(6\kappa)$, this implies the existence of $J(\epsilon)$ such that for $J \geq J(\epsilon)$, $\|u_J - u\|_\infty \leq \epsilon$. \square

A.5. Subroutine for Adding Cuts

In Table 6, we provide details regarding the subroutine for computing a new cut, which used in the R-ADP algorithm.

A.6. Some Practical Considerations in implementing R-ADP

There are two primary practical challenges that arise when implementing Algorithm 1: (1) the value function approximations are represented by an unbounded number of cuts and (2) the design of the state-sampling strategy, which becomes crucial in a high-dimensional state space.

A.6.1. Removing Redundant Cuts If we keep adding cuts to the existing approximation, some cuts become dominated by others, i.e., there exists some $k \in \{1, 2, \dots, J\}$ such that $g_k(\mathbf{y}, \boldsymbol{\gamma}) < \max_{j=1, \dots, J} g_j(\mathbf{y}, \boldsymbol{\gamma})$ for all $(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta$. It is important to remove these redundant cuts since they can lower the efficiency of solving optimization problem (20). Fortunately, the simple structure of the simplex enables us to check whether a piece is redundant efficiently and effectively. We first show how to determine whether a cut is completely dominated by another cut over the simplex.

PROPOSITION 6. $\mathbf{a}_1^T(\mathbf{y} - \mathbf{y}_1) + \mathbf{b}_1^T(\boldsymbol{\gamma} - \boldsymbol{\gamma}_1) + c_1 \geq \mathbf{a}_2^T(\mathbf{y} - \mathbf{y}_2) + \mathbf{b}_2^T(\boldsymbol{\gamma} - \boldsymbol{\gamma}_2) + c_2$ for all $(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta$ if and only if

$$c_1 - \mathbf{a}_1^T \mathbf{y}_1 - \mathbf{b}_1^T \boldsymbol{\gamma}_1 - c_2 + \mathbf{a}_2^T \mathbf{y}_2 + \mathbf{b}_2^T \boldsymbol{\gamma}_2 + \min_i \{a_{1i} - a_{2i}\}, \min_i \{b_{1i} - b_{2i}\} \geq 0.$$

Parameters: $l, c_{ij}, \rho > 0$
Data: $\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_M, \mathbf{P}_1, \mathbf{P}_2, \dots, \mathbf{P}_M$
Input: initial approximation $u_J(\mathbf{y}, \boldsymbol{\gamma}) = \max_{k=1, \dots, N_J} g_k(\mathbf{y}, \boldsymbol{\gamma})$, new sample point $(\tilde{\mathbf{y}}, \tilde{\boldsymbol{\gamma}})$
Initialization: $\mathcal{K} = \emptyset$
for $k = 1, 2, \dots, N_J$,
 if $a_{ki} - a_{kj} \leq c_{ij} \forall i, j$
 Add k into the set \mathcal{K}
 end if
end for
for $s = 1, 2, \dots, M$,
 $\bar{\mathbf{x}}_s = (\tilde{\mathbf{y}} - \mathbf{d}_s)^+ + \mathbf{P}_s^T(\tilde{\boldsymbol{\gamma}} + \min(\tilde{\mathbf{y}}, \mathbf{d}_s))$,
 $\bar{\boldsymbol{\gamma}}_s = (\tilde{\boldsymbol{\gamma}} + \min(\tilde{\mathbf{y}}, \mathbf{d}_s)) \circ (\mathbf{e} - \mathbf{P}_s \mathbf{e})$,
 Let $\mathcal{D} = \arg \max_{k=1, \dots, N_J} g_k(\bar{\mathbf{x}}_s, \bar{\boldsymbol{\gamma}}_s)$
 if $\mathcal{D} \cap \mathcal{K} = \emptyset$
 Solve $\bar{v}_J(\bar{\mathbf{x}}_s, \bar{\boldsymbol{\gamma}}_s)$ in Equation (20)
 Note the optimal solution as $\lambda_0^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*$ and v^* as the objective value
 else
 Pick any $k \in \mathcal{D} \cap \mathcal{K}$
 $\lambda_0^* = 0, \boldsymbol{\lambda}^* = \mathbf{a}_k, \boldsymbol{\mu}^* = \mathbf{e}_k, v^* = (\bar{\mathbf{x}}_s - \mathbf{y}_k)^T \mathbf{a}_k + (\bar{\boldsymbol{\gamma}}_s - \boldsymbol{\gamma}_k)^T \mathbf{b}_k + \iota_k$
 end if
 Set $\bar{c}_s = \beta_i \sum_{i=1}^n (d_{s,i} - \tilde{y}_i)^+ + \rho v^*, \bar{\mathbf{a}}_s = \lambda_0^* \mathbf{e} + \boldsymbol{\lambda}^*, \bar{\mathbf{b}}_s = \sum_{j=1}^{N_J} \mu_j^* \mathbf{b}_j$
 Set $\nabla_{\bar{\mathbf{x}}, \mathbf{y}}^s = \text{Diag}(\mathbf{1}_{\tilde{\mathbf{y}} > \mathbf{d}_s}) + \mathbf{P}_s \circ (\mathbf{1}_{\tilde{\mathbf{y}} \leq \mathbf{d}_s} \mathbf{e}^T), \nabla_{\bar{\boldsymbol{\gamma}}, \boldsymbol{\gamma}}^s = \mathbf{P}_s$
 Set $\nabla_{\bar{\boldsymbol{\gamma}}, \boldsymbol{\gamma}}^s = \text{Diag}(\mathbf{e} - \mathbf{P}_s \mathbf{e}), \nabla_{\bar{\boldsymbol{\gamma}}, \mathbf{y}}^s = \text{Diag}((\mathbf{e} - \mathbf{P}_s \mathbf{e}) \circ \mathbf{1}_{\tilde{\mathbf{y}} \leq \mathbf{d}_s})$
end for
Set $\tilde{\mathbf{a}} = \frac{1}{M} \sum_{s=1}^M (-\beta_i \sum_{i=1}^n \mathbf{1}_{\tilde{y}_i \leq d_{s,i}} \mathbf{e}_i + \rho \nabla_{\bar{\mathbf{x}}, \mathbf{y}}^s \bar{\mathbf{a}}_s + \rho \nabla_{\bar{\boldsymbol{\gamma}}, \mathbf{y}}^s \bar{\mathbf{b}}_s)$
Set $\tilde{\mathbf{b}} = \frac{1}{M} \sum_{s=1}^M (\rho \nabla_{\bar{\mathbf{x}}, \boldsymbol{\gamma}}^s \bar{\mathbf{a}}_s + \rho \nabla_{\bar{\boldsymbol{\gamma}}, \boldsymbol{\gamma}}^s \bar{\mathbf{b}}_s), \tilde{c} = \frac{1}{M} \sum_{s=1}^M \bar{c}_s$
Output: $\tilde{g}(\mathbf{y}, \boldsymbol{\gamma}) = \tilde{\mathbf{a}}^T (\mathbf{y} - \tilde{\mathbf{y}}) + \tilde{\mathbf{b}}^T (\boldsymbol{\gamma} - \tilde{\boldsymbol{\gamma}}) + \tilde{c}$

Table 6 Subroutine for Adding Cuts

Therefore, to check whether a cut is completely dominated by another cut, one just needs to perform a series of elementary operations and check one inequality, the computational effort of which is negligible compared to solving a linear program. Therefore, we can always check whether the current cut either dominates or is dominated by other cuts by checking at most $2N_J$ inequalities. Though Proposition 6 is simple to implement, it does not cover the situation where a cut is dominated by the maximum of several other cuts, which occurs more frequently. The next proposition addresses this situation.

PROPOSITION 7. $D^k \neq \emptyset$ if and only if the objective function value of the following linear program

$$\begin{aligned}
 & \min && t \\
 & \text{subject to} && \mathbf{e}^T \mathbf{y} + \mathbf{e}^T \boldsymbol{\gamma} = N, \\
 & && t \geq \mathbf{a}_l^T (\mathbf{y} - \mathbf{y}_l) + \mathbf{b}_l^T (\boldsymbol{\gamma} - \boldsymbol{\gamma}_l) + c_l - \mathbf{a}_k^T (\mathbf{y} - \mathbf{y}_k) - \mathbf{b}_k^T (\boldsymbol{\gamma} - \boldsymbol{\gamma}_k) - \iota_k, \forall l \neq k \\
 & && \mathbf{y}, \boldsymbol{\gamma} \geq 0,
 \end{aligned} \tag{43}$$

is negative.

Proof. Cut k having a dominating region in Δ means that

$$\min_{(\mathbf{y}, \boldsymbol{\gamma}) \in \Delta} \max_{l \neq k} \mathbf{a}_l^T (\mathbf{y} - \mathbf{y}_l) + \mathbf{b}_l^T (\boldsymbol{\gamma} - \boldsymbol{\gamma}_l) + c_l - \mathbf{a}_k^T (\mathbf{y} - \mathbf{y}_k) - \mathbf{b}_k^T (\boldsymbol{\gamma} - \boldsymbol{\gamma}_k) - \nu_k < 0.$$

Introducing the dummy variable t to denote the maximum, we obtain the formulation (43). \square

By solving the linear program (43) at most N_J times, one can remove all of the redundant pieces. One can perform this redundant check periodically. Our numerical implementation also employs the following strategy. We track the iterations at which a cut dominates when identifying \mathcal{D} . If a cut does not become dominating for a number of iterations greater than some threshold, we consider it potentially redundant and check for redundancy by solving (43) once. Despite these attempts to reduce the number of cuts, problem instances with more locations naturally require more cuts to accurately represent the value function $u(\cdot)$. To control the computation required for solving linear programs, we also make the practical recommendation to set an absolute upper bound on the number of cuts used in the approximation, with older cuts dropped as needed.

A.6.2. Sampling Distribution We also propose a more effective method of sampling states for the ADP algorithm beyond the naive choice of a uniform distribution over Δ . Our tests indicate that, especially when the number of locations is large, uniform sampling is unable to prioritize the sampling in important regions of the state space (for example, states with large $\boldsymbol{\gamma}$ are unlikely to occur in problem instances where the return probability is high). A reasonable strategy is to periodically simulate the ADP policy (i.e., the one implied by the current value function approximation) and collect the set of states visited under this policy — termed a *replay buffer*. On future iterations, we can sample a portion of states at which to compute cuts directly from the replay buffer. This idea is based on the notion of *experience replay* within the reinforcement learning literature (see Lin (1992) and also Mnih et al. (2015) for a recent application).

A.7. Comparative Statics

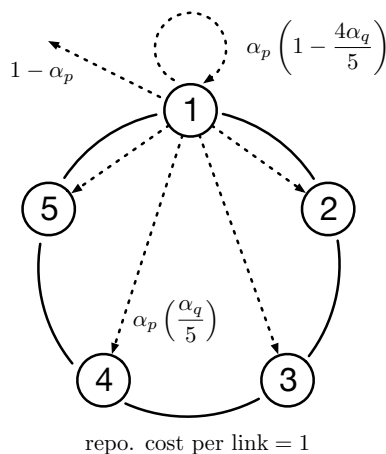


Figure 4 Network Structure used in Section A.7

Parameter Name	Value	Range
n , number of locations	5	–
N , total inventory	1.0	–
ρ , discount factor	0.95	–
β_i , lost sales cost	4.5	–
c_{ij} , repositioning cost	{1, 2}	–
α_ν , demand mean	0.3	[0.1, 0.5]
α_σ , demand volatility	1.0	[0.5, 1.5]
α_p , return fraction	0.75	[0.4, 1.0]
α_q , return uniformity	0.5	[0.0, 1.0]

Table 7 Parameter Values used in Section A.7

In this section, we aim to compare the R-ADP, myopic, and no-repositioning policies across a range of parameter settings, with the goal of studying the impacts of (1) total demand, (2) demand volatility, (3)

return fraction (i.e., fraction of vehicles returned per period), and (4) uniformity of return locations. We again use $N = 1$ and set $\rho = 0.95$. Due to the large number of MDP instances that we need to solve, we only consider $n = 5$ locations, creating a set of 9-dimensional MDPs. Let $\tilde{\nu}_i = 0.2$ for $i = 1, 2, \dots, 5$ and we set the mean demand at each location to be $\nu_i = \alpha_\nu \tilde{\nu}_i$ for some scaling parameter α_ν , so that $\sum_i \nu_i = \alpha_\nu$. Similar to before, we set $\sigma_i = \alpha_\sigma \nu_i$ for another scaling parameter α_σ . The repositioning costs c_{ij} are illustrated in Figure 4: the cost between adjacent locations is 1 and the cost between non-adjacent locations (e.g., 1 and 3 or 5 and 2) is 2. The return behavior is controlled by two parameters α_p , the fraction of rented vehicles returned (thus, $1 - \alpha_p$ is the fraction of rented vehicles that remain rented), and α_q , which we interpret as the “return uniformity.” For $\alpha_q = 1$, returns are split evenly between the 5 locations and for $\alpha_q = 0$, vehicles are returned to their rental origins. These two parameters are also illustrated in Figure 4 for location 1. We vary each of the scaling parameters α_ν , α_σ , α_p , and α_q individually; their nominal values and test ranges are summarized in Table 7. In all solved instances, we use a maximum of 1,000 cuts in the approximation while the ADP algorithm is run for 10,000 iterations¹⁷; all other algorithmic settings are as described in Section 6.3. The results are given in Figure 5. In each plot, lines illustrate absolute cost values of each policy for each parameter setting (x -axis), while the bars illustrate the percentage decrease in cost achieved by the ADP policy compared to the myopic and no-repositioning policies.

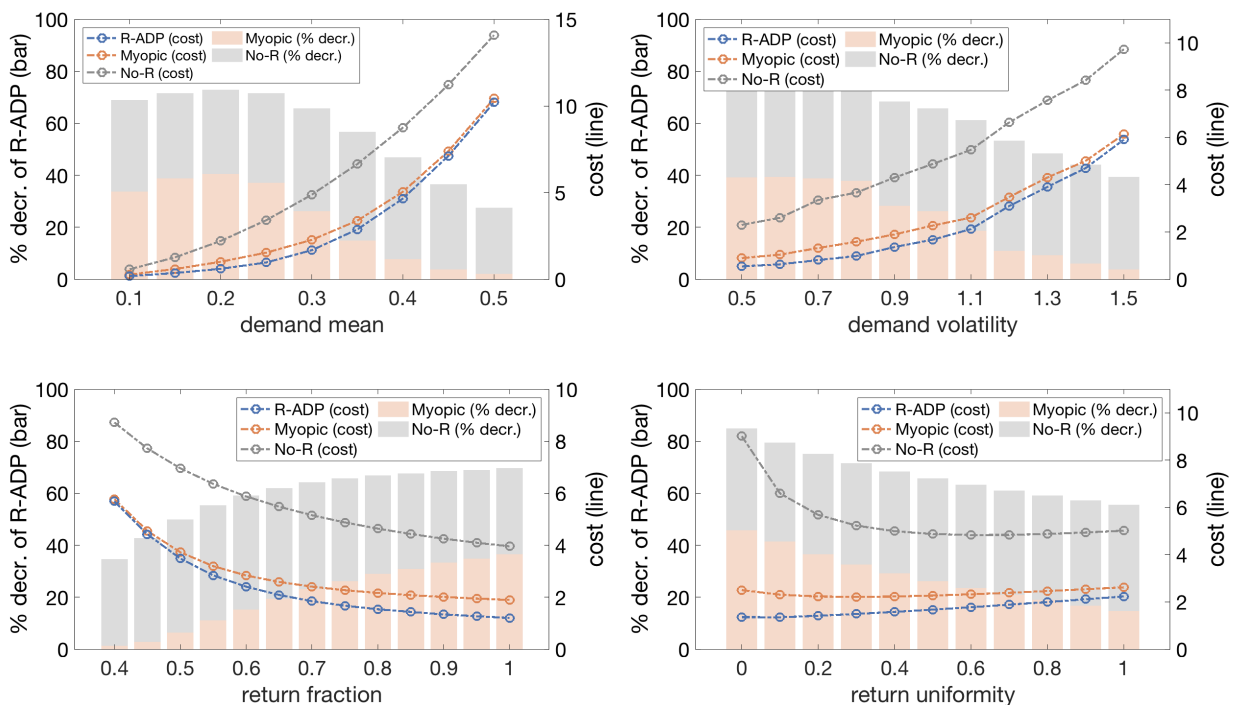


Figure 5 Impact of Parameters (Left-Axis & Bar: % Improvement of ADP; Right-Axis & Line: Raw Cost)

There are a few key takeaways from these experiments, which we now summarize. We see that when the mean demand α_ν in the system is high (greater than 45% of the total inventory), the performance of the

¹⁷ As shown in the previous section, 10,000 iterations is enough to produce a near-optimal policy for $n = 5$.

myopic policy essentially matches that of the ADP policy. This can be explained by the observation that lost sales in high demand systems are somewhat inevitable, so the impact of considering return behavior and future cost is diminished. On the other hand, when demand is between 10% and 40% of the total inventory, a substantial improvement of between 7%–40% beyond the myopic performance is observed. The largest improvement, 40%, is seen for $\alpha_\nu = 0.2$. Demand volatility α_σ has a strong influence on cost for all three policies, with the cost of the ADP policy ranging from 0.54 for $\alpha_\sigma = 0.5$ and growing tenfold to 5.90 for $\alpha_\sigma = 1.5$. We also observe that the gap between the ADP and myopic policies shrinks for higher volatility systems.

The latter two plots are related to the return behavior parameters α_p and α_q . Although the cost decreases when the return fraction α_p increases (as there is more inventory with which to satisfy demand), the improvement upon the myopic policy increases. Intuitively, given more available inventory due to fewer ongoing rentals, the ADP policy has more “opportunities” to reposition and plan for future periods. We also see that return uniformity α_q tends to increase the cost under the ADP and myopic policies, but interestingly, if no-repositioning is used, the cost is actually reduced as α_q increases in the range of $[0, 0.6]$. This can perhaps be explained by the “natural repositioning” induced by the return behaviors, an effect that disappears when active repositioning (i.e., ADP or myopic) is allowed. Similar to the case of demand volatility, the gap between the ADP and myopic policy decreases as uniformity increases.

A.8. Details of the Clustering Approach Used in the CR-ADP Algorithm

Details of the clustering approach used in Section 7.1 are provided below.

- For simplicity, suppose that n is divisible by \tilde{n} . We take a geographical clustering approach¹⁸ and group adjacent locations by setting $\mathcal{C}_k = \{(k-1)n/\tilde{n} + 1, \dots, kn/\tilde{n}\}$.
- We transform the demand by summing within clusters. For cluster k , we set $\tilde{d}_{t,k} = \sum_{i \in \mathcal{C}_k} d_{t,i}$. Moreover, we define “demand weights” χ_{ik} that conceptually represent the weight of a particular location in a cluster of locations. For each $i \in \mathcal{C}_k$, let $\chi_{ik} = \mathbf{E}d_{t,i} / \sum_{j \in \mathcal{C}_k} \mathbf{E}d_{t,j}$.
- We keep the lost sales cost unchanged: $\tilde{\beta}_i = \beta_i$. For the repositioning cost from cluster k to cluster q , we take an averaging approach:

$$\tilde{c}_{kq} = \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} \frac{1}{|\mathcal{C}_q|} \sum_{j \in \mathcal{C}_q} c_{ij}$$

- The transformed return fraction from cluster k to cluster q is given by:

$$\tilde{p}_{t,kq} = \sum_{i \in \mathcal{C}_k} \chi_{ik} \sum_{j \in \mathcal{C}_q} p_{t,ij},$$

where we first sum over \mathcal{C}_q and then take a weighted average over \mathcal{C}_k using the weights χ_{ik} defined above.

- Finally, for a location i in cluster k , we define our splitting heuristic to be: $y_i = \chi_{ik} \tilde{y}_k$.

¹⁸ Another conceivable approach is to group locations by demand (i.e., group many small locations together).

A.9. An Illustration Using Car2Go Data

In this section, we illustrate an application of the CR-ADP approach on a real world system, based on the same Car2Go dataset used in Appendix A.1. We use the same technique described in Appendix A.1 to deconstruct the GPS location time series data into individual trips that start and end within the bounds of $[-122.7, -122.6]$ longitude and $[45.5, 45.58]$ latitude. Since Car2Go is a free-floating service, we need to define discrete regions in order to apply our repositioning model and policy: to do so, we used k -means on a dataset of 64,884 GPS locations that represent either trip start or end locations to obtain $n = 200$ regions. Figure 6 shows the regions (in orange) along with the underlying GPS data (in blue).

To estimate \mathbf{P}_t , we first compute the fraction of trips started from i that end at j for all (i, j) . Next, we estimate the return fraction p_t by computing the fraction of trips that end within an hour from each region, and then average over regions to arrive at $p_t \approx 0.712$. Finally, we apply a Gaussian filter with $\sigma = 5.0$ on the resulting matrix to construct a smoother estimate. Figure 7 illustrates the estimated \mathbf{P}_t matrix using a heatmap. Examining the heatmap, we can see a prevalence of two-way trips (the bright diagonal) and one-way trips to downtown Portland (the bright vertical strip around region 50 as the destination). The mean demand at each location is simply taken to be proportional to the number of trips that started in that region; see Figure 8.

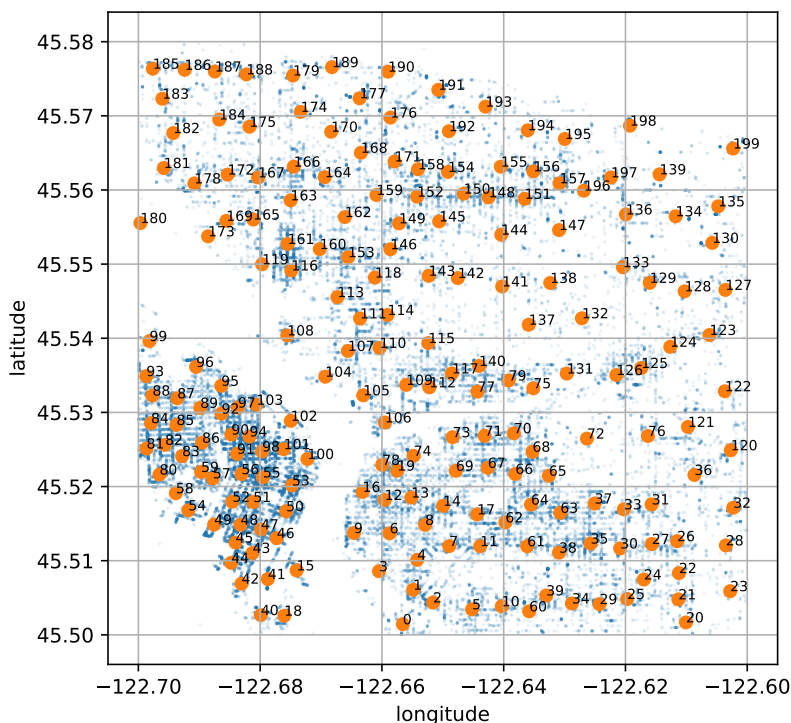


Figure 6 200 Regions (k -Means) Obtained from Car2Go GPS Data

We let the lost sales cost $\beta_i \in \{2, 4, 6, 8\}$ and the sum of mean demand $\alpha_\nu \in \{0.3, 0.4, \dots, 1.3\}$ (total inventory is normalized to be 1.0). The repositioning costs are taken to be proportional to the Euclidean

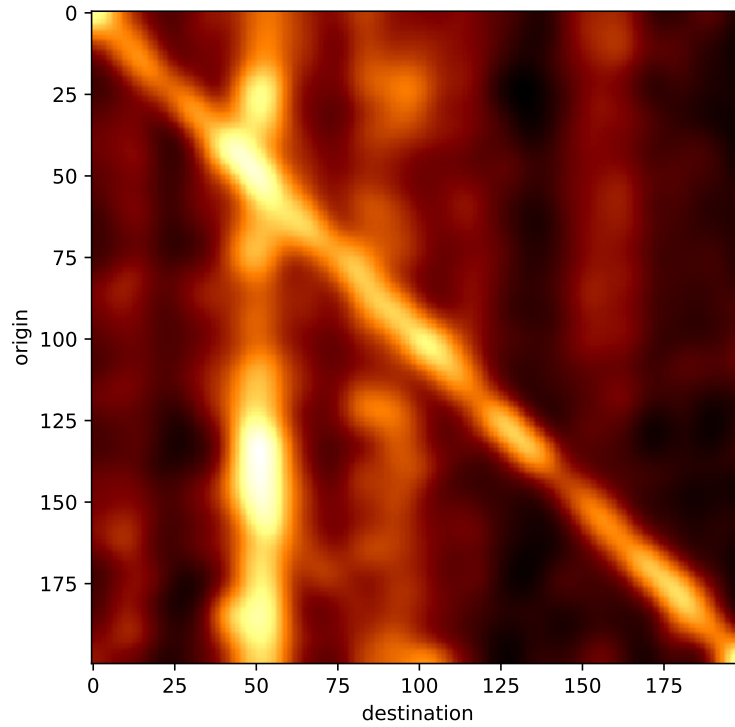


Figure 7 Heatmap of Estimated Transition P_t

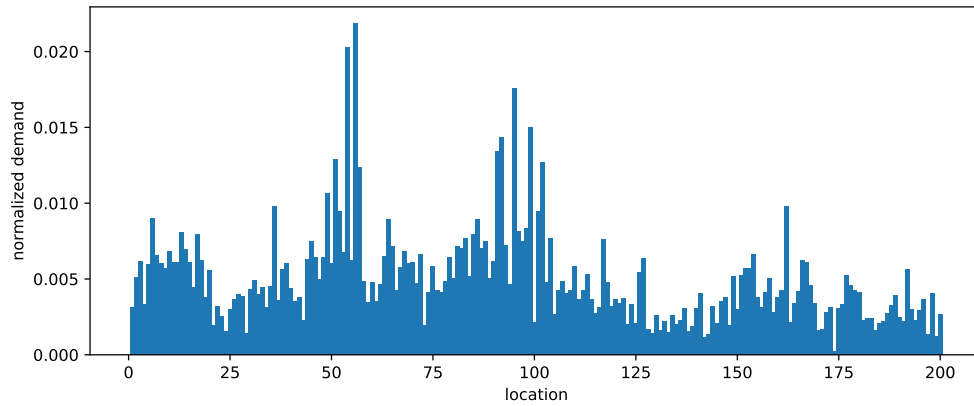


Figure 8 Normalized Demand Means

$\beta_i \mid \alpha_v$	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1	1.1	1.2	1.3
2	21.9%	36.2%	34.8%	54.0%	59.2%	67.0%	74.2%	79.3%	83.2%	86.1%	88.4%
4	18.4%	26.1%	32.7%	46.2%	62.9%	66.7%	73.7%	79.1%	83.0%	86.0%	88.2%
6	32.8%	28.7%	31.7%	47.2%	60.5%	67.4%	73.5%	79.0%	82.9%	85.9%	88.2%
8	17.9%	21.8%	33.0%	46.2%	66.2%	71.7%	73.7%	78.8%	82.8%	85.9%	88.2%

Table 8 Performance of CR-ADP as Percentage of the No-Repositioning Policy's Cost (Lower is Better)

distance between regions, ranging between 0 and 1. The discount factor is assumed to be $\rho = 0.99$. Our CR-ADP policy reduces the problem of 200 locations to a system with $\tilde{n} = 10$ clusters of equal size, with cluster k containing regions $20(k - 1) + 1$ to $20k$, as labeled in Figure 6 (these are roughly grouped by geographical location). R-ADP is run on the auxiliary MDP for 20,000 iterations and a policy for the $n = 200$ location case is then defined via the approach described in Section A.8. The results, relative to the cost of the no-repositioning policy¹⁹, are shown in Table 8. We can see that CR-ADP produces significantly lower costs in all instances. The largest relative improvements compared to no-repositioning are when the total demand is low, consistent with what is shown in Figure 5.

¹⁹The other policies are omitted due to the high computational cost of running them for $n = 200$.