# Achieving High Individual Service-Levels without Safety Stock? Optimal Rationing Policy of Pooled Resources

Jiashuo Jiang[†]    Shixin Wang[‡]    Jiawei Zhang[†]

† Department of Technology, Operations & Statistics, Stern School of Business, New York University
‡ Department of Decision Sciences and Managerial Economics, CUHK Business School, Chinese University of Hong Kong

Abstract: Resource pooling is a fundamental concept that has many applications in Operations Management for reducing and hedging uncertainty. An important problem in resource pooling is to decide (1) the capacity level of pooled resources in anticipation of random demand of multiple customers and (2) how the capacity should be allocated to fulfill customer demands after demand realization. In this paper, we present a general framework to study this two-stage problem when customers require individual and possibly different service levels. Our modeling framework generalizes and unifies many existing models in the literature, and includes second-stage allocation costs.

We propose a simple randomized rationing policy for any fixed feasible capacity level. Our main result is the optimality of this policy for very general service level constraints, including Type I and Type II constraints and beyond. The result follows from a semi-infinite linear programming formulation of the problem and its dual. As a corollary, we also prove the optimality of index policies for a large class of problems when the set of feasible fulfilled demands is a polymatroid.

## 1. Introduction

Inventory pooling is an important operational strategy that allows a firm to mitigate demand uncertainty by serving different geographic markets using a common pool of inventory. By aggregating demand across different locations, high demand from one location is likely to be offset by low demand from another. Consequently, the variability of the aggregate demand is reduced, which in turn reduces the need for safety stock (Eppen, 1979).

More specifically, consider an inventory system with $n$ locations facing independent and identically distributed (i.i.d.) demands. Each location has a target service level, say 95%, i.e., its demand must be fully satisfied with a target probability of 95%. If each location maintains a separate inventory and safety stock, the system-wide safety stock would grow *linearly* in $n$. If all locations are served by a common pool of inventory, the standard deviation of the aggregate demand is of the order $\sqrt{n}$. When the demands are i.i.d. normal distributions with a standard deviation of $\sigma$, then the total safety stock is $1.645 \cdot \sqrt{n}\sigma$ units, where the constant 1.645 is the safety factor corresponding to the target service level 95%. This so-called square-root law illustrates the benefit of inventory pooling.

In fact, the square-root law underestimates the benefit of pooling when the firm's goal is to achieve a target service level for each individual location. With a safety stock of $1.645\sqrt{n}\sigma$ units,

it is guaranteed that the aggregate demand is completely satisfied with 95% probability. When there is a system-wide shortage, it is still possible to allocate the limited inventory to locations in such a way that some locations can have their demand completely fulfilled, and those locations may achieve a service level higher than 95%. Therefore, intuitively, less safety stock is needed if the target is 95% *individual* service level! In fact, when $n = 12$ and the coefficient of variation of the normal distribution is 0.3, the system does not have to hold any safety stock to achieve 95% individual service level for all locations! The key issue here is how the inventory should be allocated when the aggregate demand exceeds the total inventory. This is the main question that we try to answer in this paper.

Inventory pooling is just one application of resource pooling in operations management. Other important applications include process flexibility (Jordan and Graves (1995), Van Mieghem (1998), Asadpour et al. (2020)), component commonality (Gerchak and Henig, 1989), transshipment (Anupindi et al., 2001), delayed differentiation (Lee, 1996), product substitution (Bassok et al., 1999); see for example Cachon and Terwiesch (2008). Individual service constraints can arise in many of these applications. This motivates us to study capacity rationing policies in a more general setting that can capture these applications.

## 1.1. Problem Formulation

We now present a general framework to model capacity allocation and demand fulfillment with *individual* service constraints. A firm serves $n$ customers, denoted by $\mathcal{N} = \{1, 2, \cdots, n\}$. The demand of customer $j \in \mathcal{N}$ is $\tilde{D}_j$ and $\tilde{\mathbf{D}} := (\tilde{D}_1, \tilde{D}_2, \ldots, \tilde{D}_n)$ follows a joint distribution $F$ with a bounded second moment. Demand of each customer can be fulfilled by utilizing one or more types of resources from the set of $m$ resources, denoted by $\mathcal{M} = \{1, 2, \cdots, m\}$.

The firm faces a two-stage decision problem. In the first stage, knowing the joint distribution $F$ but not the actual demand of the customers, the firm has to decide the capacity level of the resources $\mathbf{c} := (c_1, c_2, \ldots, c_m)$, where $c_i$ is the capacity level of resource $i \in \mathcal{M}$. The capacity investment cost is $p(\mathbf{c})$. In the second stage, the demand realizes, after which the capacity of the resources is allocated and the demand of the customers is fulfilled according to a capacity rationing policy, denoted by $\tilde{\phi}$. We denote by $s_j(\tilde{\phi}, \mathbf{c}, \mathbf{D})$ the fulfilled demand of customer $j$ under policy $\tilde{\phi}$ when the capacity level is $\mathbf{c}$ and the realized demand is $\mathbf{D} = (D_1, D_2, \ldots, D_n)$, and let $\mathbf{s}(\tilde{\phi}, \mathbf{c}, \mathbf{D}) = (s_j(\tilde{\phi}, \mathbf{c}, \mathbf{D}), j \in N)$. Notice that $s_j(\tilde{\phi}, \mathbf{c}, \mathbf{D})$ can be a random variable even for fixed demand $\mathbf{D}$ if $\tilde{\phi}$ is allowed to be a randomized policy. Similarly, we denote by $y_{ij}(\tilde{\phi}, \mathbf{c}, \mathbf{D})$ the allocation of resource $i$ to customer $j$ and let $\mathbf{y}(\tilde{\phi}, \mathbf{c}, \mathbf{D}) = (y_{ij}(\tilde{\phi}, \mathbf{c}, \mathbf{D}), i \in \mathcal{M}, j \in \mathcal{N})$. The allocation cost is denoted by $f(\mathbf{y}(\tilde{\phi}, \mathbf{c}, \mathbf{D}))$. We assume that $f(\cdot)$ is a linear function throughout the paper.

Regardless of the rationing policy used, the allocation must satisfy the following constraints. More specifically, given $\mathbf{c}$ and $\mathbf{D}$, the set of all feasible fulfilled demands and resource allocation is denoted by $P(\mathbf{c}, \mathbf{D})$. By specializing the choices $P(\mathbf{c}, \mathbf{D})$, the feasible set captures the capacity consumption and demand fulfillment constraints of many capacity allocation models such as inventory pooling, process flexibility, and assemble-to-order, etc., as illustrated below.

- In inventory pooling, $\mathcal{N}$ is the set of locations, $\mathcal{M}$ is a singleton, and

$$P(\mathbf{c}, \mathbf{D}) = \left\{ (\mathbf{s}, \mathbf{y}) \geq 0 : \quad \sum_{j=1}^{n} y_j \leq c, \quad \mathbf{y} = \mathbf{s} \quad \mathbf{s} \leq \mathbf{D} \right\}. \tag{1}$$

  This special case will also be referred to as the single resource allocation problem.
- In process flexibility, $\mathcal{N}$ is the set of products, $\mathcal{M}$ is the set of plants, and

$$P(\mathbf{c}, \mathbf{D}) = \left\{ (\mathbf{s}, \mathbf{y}) \geq 0 : \quad \mathbf{s} \leq \mathbf{D}, \quad \sum_{j \in \mathcal{N}:(i,j) \in E} y_{ij} \leq c_i \quad \forall i \in \mathcal{M}, \quad \sum_{i \in \mathcal{M}:(i,j) \in E} y_{ij} = s_j \quad \forall j \in \mathcal{N} \right\} \tag{2}$$

  where the set $E$ represents the design of the flexible system: $(i,j) \in E$ if product $j$ can be produced by plant $i$.
- In an assemble-to-order system, $\mathcal{N}$ is the set of end products, $\mathcal{M}$ is the set of components,

$$P(\mathbf{c}, \mathbf{D}) := \{ (\mathbf{s}, \mathbf{y}) \geq 0 : \quad A\mathbf{s} \leq \mathbf{c}, \quad \mathbf{s} \leq \mathbf{D}, \quad y_{ij} = A_{ij} s_j, \quad \forall i \in \mathcal{M}, j \in \mathcal{N} \} \tag{3}$$

  where $A_{ij} \geq 0$ is the amount of component $i$ that each unit of product $j$ requires. In a special case, the so-called generalized W-system, $A$ is specialized as

$$A = \begin{bmatrix} I_{n \times n} \\ \mathbf{1}_n^T \end{bmatrix} \tag{4}$$

  where $\mathbf{1}_n$ is the $n$-dimensional column vector of all ones and $I_{n \times n}$ is the $n \times n$ identity matrix. In this system, each end product $j$ requires two components, a product-specific component $j$ and the component $n+1$, the latter of which is common to all end products.

Clearly, not all demands can always be fulfilled, but the firm is obligated to achieve a target individual service level $\beta_j \in (0, 1)$ for each customer $j \in \mathcal{N}$. This service level constraint can be formally formulated as

$$\mathrm{E}_{\tilde{\boldsymbol{\phi}}, \tilde{\mathbf{D}}}[R_j(s_j(\tilde{\boldsymbol{\phi}}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \geq \beta_j \tag{5}$$

where $R_j : \mathcal{R}_+^2 \to \mathcal{R}_+$ is called the service measure function of customer $j$. Constraint (5) unifies different types of service level constraints in the operations management literature. For example, when $R_j(s_j, D_j) = 1_{s_j \geq D_j}$, constraint (5) defines the so-called Type I service level constraint, i.e., the demand of customer $j$ must be completely satisfied with probability at least $\beta_j$. Similarly,

Type II service level can be defined by letting $R_j(s_j, D_j) = s_j / \mathrm{E}[\tilde{D}_j]$, which measures the fraction of the expected demand that can be satisfied. In contrast, choosing $R_j(s_j, D_j) = s_j / D_j$ allows us to measure the fraction of the actual demand that can be fulfilled, which we name as Type III service level constraint. It is straightforward to verify that these functions all satisfy the following conditions.

ASSUMPTION 1.

  a. *For any* $j \in \mathcal{N}$, $R_j(s_j, D_j)$ *is non-decreasing and upper semi-continuous in* $s_j$, *for any fixed* $D_j$.

  b. *For any* $\mathbf{c}$ *and* $\mathbf{D}$, $P(\mathbf{c}, \mathbf{D})$ *is a compact set.*

The firm's problem is to decide capacity level $\mathbf{c}$ and rationing policy $\tilde{\phi}$ to minimize the first stage capacity investment cost $p(\mathbf{c})$ and the expected second stage allocation cost subject to the individual service constraints, which can be formulated as

$$\inf_{\mathbf{c} \geq 0, \tilde{\phi}} \quad p(\mathbf{c}) + \mathrm{E}_{\tilde{\phi}, \tilde{\mathbf{D}}}[f(\mathbf{y}(\tilde{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] \tag{6}$$

$$\text{s.t.} \quad \mathrm{E}_{\tilde{\phi}, \tilde{\mathbf{D}}}[R_j(s_j(\tilde{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \geq \beta_j \quad \forall j \in \mathcal{N} \tag{6a}$$

$$(s_j(\tilde{\phi}, \mathbf{c}, \mathbf{D}), \mathbf{y}(\tilde{\phi}, \mathbf{c}, \mathbf{D})) \in P(\mathbf{c}, \mathbf{D}) \quad \forall \mathbf{D}.$$

The formulation should also specify the set of feasible (randomized) policies, which will be discussed in Section 2.

Although (6) is formulated as a single-period model, it is possible to approximate it by a periodic-review infinite time horizon problem as follows. Assume that the capacity is perishable, i.e., unused capacity in the previous period can not be used to satisfy future demands, and unmet demands are lost. Denote the demand in period $t$ by $\tilde{\mathbf{D}}^{(t)}$ and assume $\tilde{\mathbf{D}}^{(1)}, \tilde{\mathbf{D}}^{(2)}, \cdots, \tilde{\mathbf{D}}^{(t)}, \cdots$, are i.i.d random variables. The fulfilled demand of customer $j$ and the allocation in period $t$ is denoted by

$$s_j^{(t)}(\mathbf{c}, \mathbf{D}^{(1:t)}, \mathbf{s}^{(1:t-1)}, \mathbf{y}^{(1:t-1)}) \quad \text{and} \quad \mathbf{y}^{(t)}(\mathbf{c}, \mathbf{D}^{(1:t)}, \mathbf{s}^{(1:t-1)}, \mathbf{y}^{(1:t-1)})$$

where $\mathbf{D}^{(1:t)} = (\mathbf{D}^{(1)}, \cdots, \mathbf{D}^{(t)})$, $\mathbf{s}^{(1:t-1)} = (\mathbf{s}^{(1)}, \cdots, \mathbf{s}^{(t-1)})$ and $\mathbf{y}^{(1:t-1)} = (\mathbf{y}^{(1)}, \cdots, \mathbf{y}^{(t-1)})$, which shows explicitly that the demand fulfillment and resource allocation decisions in period $t$ can depend on realized demands up to time $t$ and on previous fulfillment and resource allocation decisions up to time $t-1$. With these notations, formulation (6) can be approximated by

$$\inf_{\mathbf{c} \geq 0, \mathbf{s} \geq 0, \mathbf{y} \geq 0} \quad p(\mathbf{c}) + \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} f(\mathbf{y}^{(t)}(\mathbf{c}, \mathbf{D}^{(1:t)}, \mathbf{s}^{(1:t-1)}, \mathbf{y}^{(1:t-1)})) \tag{7}$$

$$\text{s.t.} \quad \liminf_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} R_j(s_j^{(t)}(\mathbf{c}, \mathbf{D}^{(1:t)}, \mathbf{s}^{(1:t-1)}, \mathbf{y}^{(1:t-1)}), D_j^{(t)}) \geq \beta_j, \quad \forall j \in \mathcal{N} \tag{7a}$$

Note that the service level constraint (7a) is defined in an asymptotic sense, which implies the following definition of *asymptotic* feasibility for the single-period model.

DEFINITION 1. A capacity level **c** is *asymptotically* feasible as long as for any $\epsilon > 0$, there exists a rationing policy $\tilde{\phi}_\epsilon$ such that

$$\mathrm{E}_{\tilde{\phi}_\epsilon, \tilde{\mathbf{D}}}[R_j(s_j(\tilde{\phi}_\epsilon, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \geq \beta_j - \epsilon, \quad \forall j \in \mathcal{N} \tag{8}$$

Obviously, asymptotic feasibility does not immediately imply feasibility, i.e., constraint (6a). However, we shall prove that under an additional yet mild assumption stated in the next section, asymptotic feasibility is equivalent to feasibility, i.e., constraint (8) is equivalent to constraint (6a). To the best of our knowledge, this is the first time that such an equivalence for the service level constraint is formally proved in the literature.

## 1.2. Previously Known Results and Closely Related Literature

There are different models concerning service level constraints in the literature, which can be captured as special cases of our general framework. To the best of our knowledge, existing literature on individual service level constraints considers only first stage capacity investment cost, but not the second stage resource allocation cost.

Baker (1985) was among the first to discuss the impact of individual Type I service constraints on safety stock level in the context of an assemble to order system. A more formal analysis appeared in Baker et al. (1986) that considers a W-system with two products and three components. Therefore, capacity rationing only matters for the common component. However, optimal rationing policy is not discussed in the paper. In their analysis of safety stock level, it is assumed that the priority is always given to the product with lower realized demand. An optimal allocation policy is derived for the same problem, i.e., the W-system with two products and three components, in Gerchak et al. (1988) when the target service levels of the two products are equal, and in Mirchandani and Mishra (2002) when the target service levels of the two products can be different.

Individual Type I service constraints have also been addressed in inventory pooling. A chance-constrained stochastic program was formulated in Swaminathan and Srinivasan (1999). However, the number of decision variables grows exponentially in the number of customers. Computational results are reported only when the number of customers is small, i.e., two or three. A closed-form expression for the optimal inventory level is derived when the number of customers is two and demands are i.i.d uniform distributions. Alptekinoğlu et al. (2013) prove that priority policies are optimal under which customers are served according to a priority list. An optimal policy is derived when demands are i.i.d. They also compare the performance of anticipative vs responsive priority policies. We will define and discuss these two types of priority policies in detail in Section 3.

Much progress has recently been made when individual service constraints are of Type II. Hou et al. (2009) study the single-resource allocation problem in the context of wireless networks problem with quality of service (QoS) constraints, which is the same as the infinite horizon multi-period

formulation of the Type II service constraints. They propose the so-called largest-debt-first policy and prove its optimality for the single-resource allocation problem. Their analysis is based on a novel application of the celebrated Blackwell's approachability theorem (Blackwell, 1956). Zhong et al. (2017) apply a similar approach and analyze the safety stock level in inventory pooling with individual Type II service constraints. They show that a randomized anticipative priority policy is optimal in this setting. The approach and results are extended, in a highly non-trivial way, by Lyu et al. (2019) to study capacity rationing policy in process flexibility.

There is a stream of literature in inventory management that studies the problem of fulfilling the demand of multiple customers with individual Type I service guarantees; see for example Agrawal and Cohen (2001) and Zhang (1997). However, most of these papers focus on inventory optimization under given inventory rationing policies and these policies are not necessarily optimal. For example, Agrawal and Cohen (2001) propose a heuristic policy, called the fair-share allocation policy, while Zhang (1997) assumes that customer demands are fulfilled according to a pre-specified priority list.

If we require in our model that the entire demand of all customers must be met with a given probability, capacity rationing is no longer needed and our problem can then be formulated as a two-stage joint chance-constrained stochastic program, which has been studied in Gurvich et al. (2010) and Liu et al. (2016).

### 1.3.  Our Results

We propose in this paper a simple rationing policy, called the Max-Weighted-Service policy, for formulation (6). Our policy assigns a random weight to each customer and based on the weights we solve, after demand realization, a deterministic capacity allocation and demand fulfillment problem to maximize a weighted service measure function. The random weight is sampled from a sufficiently large set that can be constructed offline.

Our main result is to show the Max-Weighted-Service policy is asymptotically optimal for a very general class of capacity allocation and demand fulfillment problems with individual service constraints. The generality of our model formulation is similar to the so-called newsvendor networks models proposed by Mieghem and Rudi (2002) (see also Bassamboo et al. (2010)). Indeed, our model is applicable in inventory pooling, process flexibility, assemble to order, transshipment, substitution, etc, as long as the fulfilled demand can be modeled as a linear transformation of capacity, i.e., the feasible set $P(\mathbf{c}, \mathbf{D})$ is a bounded polyhedron.

Unlike the models of Mieghem and Rudi (2002) that penalize unsatisfied demand in the cost function, our model explicitly imposes individual service constraint for each customer. And the

---

[0] After presenting this work at NYU Shanghai on April 19th, 2019, it was brought to our attention by Professor Renyu Zhang that Blackwell's approachability theorem has also been applied to study single-resource pooling with Type I service constraints by Lyu et al. (2017).

service constraints can be defined in a variety of ways. Indeed, our Max-Weighted-Service policy is asymptotically optimal under very mild conditions on the service measure functions, which are satisfied by both Type I and Type II service levels. As discussed in the previous subsection, Type I and Type II service levels have usually been analyzed separately in the literature. Our approach allows a unified treatment for both service levels, and beyond. Besides these two metrics that are commonly used in practice and studied in the literature, we also allow the service level of a customer to depend on, for example, the probability that its demand is fully satisfied as well as the probability that a certain fraction of its demand is satisfied.

Despite the generality of the model, our approach to derive the policy is simple. We formulate the problem of finding an optimal randomized policy, for a fixed capacity level, as a semi-infinite linear program. The decision variable can be interpreted as the probability measure over the set of all possible deterministic policies, not just all possible priority lists. (Priority policies are not always optimal for our general model.) Although this formulation is natural, it appears to be new in the literature that addresses individual service constraints. Randomized policies have been studied for various special cases of our model, see for example Swaminathan and Srinivasan (1999), Alptekinoğlu et al. (2013), Zhong et al. (2017), Lyu et al. (2019), but their formulations are different than ours. For example, for inventory pooling with Type I service constraints, Swaminathan and Srinivasan (1999) partition the support of demand into different regions, and the decision variable is, for each demand region, the probability of choosing a particular priority list. Alptekinoğlu et al. (2013)) takes a similar approach. We discuss the difference between our approach and those of Zhong et al. (2017) and Lyu et al. (2019) below.

As a corollary of our main result, we show that randomized anticipative (responsive, respectively) index policies are asymptotically optimal when all individual service constraints are of Type II (Type III, respectively) and when the feasible set $P(\mathbf{c}, \mathbf{D})$ can be characterized by a polymatroid. Even with this additional assumption on $P(\mathbf{c}, \mathbf{D})$, our model still captures a wide range of problems as special cases such as inventory pooling, process flexibility, commonality in the generalized W-system in assemble-to-order, capacity planning of more general network, etc. For the special cases of single-resource pooling and process flexibility without second-stage allocation costs, our policy recovers those in Hou et al. (2009), Zhong et al. (2017), and Lyu et al. (2019). However, our policy is derived using a different approach. For example, Lyu et al. (2019) uses the infinite-time horizon model to study the single-period model. They derive an allocation policy for the infinite-time horizon model. They then use their policy to derive a sufficient condition for a given capacity being feasible for the infinite-time horizon model, which corresponds to the notion of asymptotically feasible defined in our paper for the single-period model. Their analysis appears to be specific to the process flexibility problem and is much more involved than the analysis of

the single-resource allocation problem in Zhong et al. (2017). In contrast, our semi-infinite linear programming formulation and its dual allow us to derive the necessary and sufficient condition of a given capacity level for the single-period problem directly. Our allocation policy is motivated by applying the stochastic gradient descent (SGD) algorithm to the dual problem. Moreover, the SGD-based approach proves the asymptotic optimality of the policy for a much more general model that even includes second-stage allocation costs, while Blackwell's approachability theorem focuses on finding a feasible path to approach the target set without concerning optimality of this path with regard to the allocation cost. Also, the connection between the dual variable and the optimal allocation policy appears to be new.

Based on our duality result, we develop a minimax stochastic programming formulation and apply an existing first-order algorithm to compute an optimal or near-optimal capacity level. Numerical results show that the algorithm converges to a globally optimal solution whenever the objective function is convex, as predicted by existing theory. When the objective function is non-convex, we propose heuristics to compute near-optimal solutions.

Our framework is general and can be applied to different problems and different service level constraints. Though we have different assumptions in the following sections, they are all mild. To put it succinctly, the sufficient and necessary condition for a given capacity to be asymptotically feasible (Theorem 1) or feasible (Theorem 2), the asymptotic optimality of the Max-Weighted-Service policy (Theorem 3) and the minimax formulation for solving the optimal capacity (Theorem 5) apply to settings including inventory pooling, flexible production and the assemble-to-order problems for service levels including Type I, Type II and Type III. The optimality of an index policy (in Theorem 4) requires a stronger assumption and applies to inventory pooling, flexible production problems and the generalized W-system ATO problems under Type II or Type III service levels, and does not apply to general ATO problems or Type I service level.

## 2.  Randomized Rationing Policy and Problem Reformulation

We begin this section by formally formulating the set of feasible rationing policies in formulation (6). A *deterministic* policy is a function $\phi$ from $\mathbf{R}_+^{m+n}$ to $\mathbf{R}_+^n \times \mathbf{R}_+^{nm}$ such that for any capacity level $\mathbf{c}$ and any realized demand $\mathbf{D}$,

$$\phi(\mathbf{c}, \mathbf{D}) = (\mathbf{s}(\phi, \mathbf{c}, \mathbf{D}), \mathbf{y}(\phi, \mathbf{c}, \mathbf{D})) \in P(\mathbf{c}, \mathbf{D}).$$

We also denote by $\mathbf{s}(\phi, \mathbf{c}, \mathbf{D})$ the demand fulfillment and $\mathbf{y}(\phi, \mathbf{c}, \mathbf{D})$ the allocation under the deterministic policy $\phi$ for fixed $\mathbf{c}$ and $\mathbf{D}$. We denote the set of all deterministic policies by $\Phi$.

A *randomized* policy is determined by a probability measure $\lambda$ over $\Phi$ such that any (measurable) subset of deterministic policies $\hat{\Phi} \subseteq \Phi$ is chosen with probability $\lambda(\hat{\Phi})$. Such a randomized policy is

denoted by $\tilde{\phi}_\lambda$ or simply $\lambda$. Therefore, optimization over randomized policies can be reformulated as an optimization over probability measures.

Under a deterministic policy $\phi \in \Phi$, the service level of customer $j$ is given by $E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)]$. Therefore, under a randomized policy $\tilde{\phi}_\lambda$, the service level of customer $j$ is

$$E_{\tilde{\phi}_\lambda, \tilde{\mathbf{D}}}[R_j(s_j(\tilde{\phi}_\lambda, \mathbf{c}, \tilde{\mathbf{D}}), D_j)] = \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\phi). \tag{9}$$

It follows that problem (6) can be reformulated as

$$\inf_{\mathbf{c} \geq 0, \lambda \in \chi} \quad p(\mathbf{c}) + \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda(\phi) \tag{10}$$

$$\text{s.t.} \quad \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\phi) \geq \beta_j \quad \forall j \in \mathcal{N} \tag{10a}$$

where $\chi = \{\lambda \geq 0 : \int_{\phi \in \Phi} d\lambda(\phi) = 1\}$. Notice that constraint (10a) appears to be bilinear in $(E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)], d\lambda)$, which is non-convex.

In the remainder of this section, we assume that capacity $\mathbf{c}$ is fixed and establish conditions for checking feasibility of a given capacity level. We also present in Section 3 an optimal rationing policy when the given capacity level is feasible. These results will be used in Section 4 for the computation of optimal or near-optimal capacity levels. The results can be obtained by considering the following semi-infinite linear program

$$\inf_{\lambda \in \chi} \quad \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda(\phi) \tag{11}$$

$$\text{s.t.} \quad \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\phi) \geq \beta_j, \quad \forall j \in \mathcal{N} \tag{11a}$$

and its Lagrangian dual formulation. Introducing the Lagrangian dual multipliers $w_j$ for constraints (11a) for each $j \in \mathcal{N}$, we obtain the Lagrangian dual formulation of (11)

$$\sup_{\mathbf{w} \geq 0} \quad \inf_{\lambda \in \chi} L(\mathbf{w}, \lambda) \tag{12}$$

where

$$L(\mathbf{w}, \lambda) := \sum_{j \in \mathcal{N}} w_j \cdot \beta_j + \int_{\phi \in \Phi} F(\mathbf{w}, \phi) d\lambda(\phi) \tag{13}$$

denotes the Lagrangian dual function and

$$F(\mathbf{w}, \phi) := E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] - \sum_{j \in \mathcal{N}} w_j \cdot E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)]. \tag{14}$$

It is clear that for fixed $\mathbf{w}$, $\inf_{\lambda \in \chi} L(\mathbf{w}, \lambda) = \sum_{j \in \mathcal{N}} w_j \cdot \beta_j + \inf_{\phi \in \Phi} F(\mathbf{w}, \phi)$ and thus the dual problem can be reformulated as

$$\sup_{\mathbf{w} \geq 0} \inf_{\lambda \in \chi} L(\mathbf{w}, \lambda) = \sup_{\mathbf{w} \geq 0} \sum_{j \in \mathcal{N}} w_j \cdot \beta_j + \inf_{\phi \in \Phi} F(\mathbf{w}, \phi) \tag{15}$$

Indeed, the above dual formulation can be further simplified. To that end, We first define a deterministic optimization problem which we call the Max-Weighted-Service problem. Specifically, for any given $\mathbf{w} \geq 0, \mathbf{c}$, and $\mathbf{D}$, define

$$g(\mathbf{w}, \mathbf{c}; \mathbf{D}) = \min_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \; f(\mathbf{y}) - \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j) \tag{16}$$

Under Assumption 1, problem (16) always attains its minimum in the compact set $P(\mathbf{c}, \mathbf{D})$. For any $\mathbf{w} \geq 0$, we define a deterministic policy $\boldsymbol{\phi}_{\mathbf{w}}$ such that $\boldsymbol{\phi}_{\mathbf{w}}(\mathbf{c}, \mathbf{D}) = (\mathbf{s}_{\mathbf{w}}^*(\mathbf{c}, \mathbf{D}), \mathbf{y}_{\mathbf{w}}^*(\mathbf{c}, \mathbf{D}))$ for any $\mathbf{c}$ and $\mathbf{D}$, where $(\mathbf{s}_{\mathbf{w}}^*(\mathbf{c}, \mathbf{D}), \mathbf{y}_{\mathbf{w}}^*(\mathbf{c}, \mathbf{D}))$ denotes an optimal solution of (16). (When (16) has multiple optimal solutions, ties are broken arbitrarily so that $(\mathbf{s}_{\mathbf{w}}^*(\mathbf{c}, \mathbf{D}), \mathbf{y}_{\mathbf{w}}^*(\mathbf{c}, \mathbf{D}))$ is uniquely defined.)

We show in the next lemma that it suffices to focus on the Max-Weighted-Service problem to solve the dual problem (15), where the proof is relegated to Appendix A.

LEMMA 1. *For any fixed* $\mathbf{w} \geq 0$, *it holds that*

$$\inf_{\boldsymbol{\phi} \in \Phi} \; F(\mathbf{w}, \boldsymbol{\phi}) = \; F(\mathbf{w}, \boldsymbol{\phi}_{\mathbf{w}}) = E_{\tilde{\mathbf{D}}}[g(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})]. \tag{17}$$

In fact, when considering the feasibility of a given capacity level, we can further simplify (12) by assuming zero allocation cost. We are now ready to present our first result regarding the asymptotic feasibility of a given capacity level $\mathbf{c}$.

THEOREM 1. *Under Assumption 1, a given capacity level* $\mathbf{c}$ *is asymptotically feasible if and only if*

$$E_{\tilde{\mathbf{D}}}\big[ \max_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j)\big] \geq \sum_{j \in \mathcal{N}} w_j \beta_j \quad \text{for all} \;\; \mathbf{w} \geq 0 \tag{18}$$

*Proof:* When a capacity level $\mathbf{c}$ is asymptotically feasible, from (8), it is clear that we have

$$\inf_{\lambda \in \chi} \; \sum_{j \in \mathcal{N}} w_j \cdot \left( \beta_j - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\boldsymbol{\phi}) \right) \leq 0$$

for each fixed $\mathbf{w} \geq 0$. Thus, it holds that

$$\sup_{\mathbf{w} \geq 0} \; \inf_{\lambda \in \chi} \; \sum_{j \in \mathcal{N}} w_j \cdot \left( \beta_j - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\boldsymbol{\phi}) \right) \leq 0 \tag{19}$$

From Lemma 1 (assuming zero allocation cost), we immediately tell that (18) holds.

We now prove the reverse direction. If (18) holds, then (19) holds from Lemma 1. We define the set $W = \{\mathbf{w} \geq 0 : \sum_{j \in \mathcal{N}} w_j \leq 1\}$. Clearly, we have that

$$\max_{\mathbf{w} \in W} \; \inf_{\lambda \in \chi} \; \sum_{j \in \mathcal{N}} w_j \cdot \left( \beta_j - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\boldsymbol{\phi}) \right) \leq 0 \tag{20}$$

Obviously, the set $\chi$ is a convex set. Moreover, note that $W$ is a convex compact set and the objective function in the above problem is linear in $\mathbf{w}$ (resp. $\lambda$) when $\lambda$ is fixed. Then by Sion's

minimax theorem (Sion, 1958), we can interchange the order of max and inf on the left-hand side of (20). Thus, we have

$$\inf_{\lambda \in \chi} \max_{\mathbf{w} \in W} \sum_{j \in \mathcal{N}} w_j \cdot \left( \beta_j - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda(\boldsymbol{\phi}) \right) \leq 0$$

For any $\epsilon > 0$, there exists a randomized policy $\lambda_\epsilon \in \chi$ such that

$$\sup_{\mathbf{w} \in W} \sum_{j \in \mathcal{N}} w_j \cdot \left( \beta_j - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_\epsilon(\boldsymbol{\phi}) \right) \leq \epsilon \tag{21}$$

We now claim that $\lambda_\epsilon$ achieves a service level at least $\beta_j - \epsilon$ for each $j \in \mathcal{N}$, i.e.

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_\epsilon(\boldsymbol{\phi}) \geq \beta_j - \epsilon. \tag{22}$$

Otherwise, suppose there exists a $j' \in \mathcal{N}$ such that (22) does not hold. Then we define $\hat{\mathbf{w}} \in W$ such that $\hat{w}_{j'} = 1$ and $\hat{w}_j = 0$ for all $j \neq j'$. It is clear that

$$\epsilon < \hat{w}_{j'} \cdot \left( \beta_{j'} - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_{j'}(s_{j'}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_{j'})] d\lambda_\epsilon(\boldsymbol{\phi}) \right)$$

$$\leq \sup_{\mathbf{w} \in W_\epsilon} \sum_{j \in \mathcal{N}} w_j \cdot \left( \beta_j - \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_\epsilon(\boldsymbol{\phi}) \right) \leq \epsilon$$

where the last inequality follows from (21). This is a contradiction. Thus, we prove that $\mathbf{c}$ is asymptotically feasible, which completes our proof. $\square$

Recall that asymptotic feasibility is weaker than feasibility. In what follows, we present sufficient conditions for feasibility. To that end, we prove strong duality between the primal-dual pair (11) and (12) under the following assumption.

ASSUMPTION 2.

a. *For any $j \in \mathcal{N}$, the service measure function $R_j$ satisfies:*
   - $R_j(s_j, D_j)$ *is non-decreasing in $s_j$, for any fixed $D_j$;*
   - *there exists a finite set of parameters $0 = a_{j,1} < a_{j,2} < \cdots < a_{j,K_j} = 1$ such that for any fixed $D_j$, $R_j(s_j, D_j)$ is linear in $s_j$ when $s_j \in [a_{j,l}D_j, a_{j,l+1}D_j)$ for any $l \in \{1, 2, \ldots, K_j - 1\}$;*
   - *there exists a constant $C_1 > 0$ such that $R_j(s_j, D_j) \leq C_1 \cdot \max\{1, D_j\}$ for any $\mathbf{D}$ and any $s_j \leq D_j$.*

b. *$P(\mathbf{c}, \mathbf{D})$ is a bounded polyhedron of $(\mathbf{s}, \mathbf{y})$ defined by a set of linear inequalities on $(\mathbf{s}, \mathbf{y}, \mathbf{c}, \mathbf{D})$, including the constraints $\mathbf{s} \leq \mathbf{D}$ and $(\mathbf{s}, \mathbf{y}) \geq 0$, and there exists a constant $C_2 > 0$ such that $\|(\mathbf{s}, \mathbf{y})\|_2^2 \leq C_2 \cdot \|\mathbf{D}\|_2^2$ for any $\mathbf{D}$ and any $(\mathbf{s}, \mathbf{y}) \in P(\mathbf{c}, \mathbf{D})$.*

Assumption 2a requires $R_j(s_j, D_j)$ to be piece-wise linear in $s_j$ for any fixed $D_j$. The breakpoints are defined based on the ratio $s_j/D_j$, which denotes the fraction of the fulfilled demand of customer $j$. It is satisfied by Type I, Type II, Type III service measure functions. Moreover, models such as inventory pooling, process flexibility and assemble-to-order all satisfy Assumption 2b.

THEOREM 2. *Under Assumption 2, strong duality holds between* (11) *and* (12) *for any given capacity level* $\mathbf{c} \geq 0$. *Specifically,* (11) *is feasible if and only if the objective value of* (12) *is finite, and when* (11) *is feasible, the objective values of* (11) *and* (12) *are the same.*

The proof is relegated to Appendix B. We note that strong duality for semi-infinite linear programming under various conditions has been studied in the literature; see for example Shapiro (2001) and Martin et al. (2016). However, we have not found a simple way to verify these conditions. For example, in order to apply the strong duality results of Shapiro (2001), we have to show certain closedness or compactness properties of a topological space on the set of deterministic policies $\Phi$. Our proof essentially shows certain compactness of a related set. However, we choose to apply Sion's minimax theorem (Sion, 1958) to avoid introducing additional concepts required by existing conditions on the semi-infinite linear programming duality and our proof appears to be slightly simpler.

By Theorem 2, problem (11) is feasible if and only if the objective value of (12) is finite. In fact, we can simplify the condition by assuming zero allocation cost in (11). The following is an immediate corollary of Theorem 2 and Lemma 1.

COROLLARY 1. *Under Assumption 2, a given capacity level* $\mathbf{c}$ *is feasible if and only if*

$$E_{\tilde{\mathbf{D}}}\big[ \max_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j) \big] \geq \sum_{j \in \mathcal{N}} w_j \beta_j \quad \text{for all} \quad \mathbf{w} \geq 0 \tag{23}$$

Corollary 1 can be used to develop numerical procedures to check whether or not a given capacity level is feasible. In fact, for fixed $\mathbf{c}$, condition (23) is equivalent to

$$\max_{\mathbf{w} \geq 0} \sum_{j \in \mathcal{N}} w_j \beta_j - E_{\tilde{\mathbf{D}}}\big[ \max_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j) \big] \leq 0 \tag{24}$$

It is clear that $\max_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j)$ is convex in $\mathbf{w}$ for any fixed $\mathbf{c}$ and $\mathbf{D}$. Thus, checking the feasibility of a fixed $\mathbf{c}$ amounts to solving a concave maximization problem. In Appendix D, we further illustrate how Corollary 1 can be used to find minimum capacity for problems studied in Mirchandani and Mishra (2002) and Swaminathan and Srinivasan (1999).

Although Corollary 1 enables us to check whether or not a given capacity level is feasible, it does not explicitly guide *how* the service guarantees can be achieved. Indeed, the result was derived from an existence proof by the strong duality between (11) and (12), and it does not immediately suggest an allocation policy. However, as we shall discuss in the next section, the dual formulation (12) is instrumental for us to derive a capacity rationing policy for a feasible capacity level.

## 3. Max-Weighted-Service Policy

In this section, we solve problem (11) for a fixed and feasible $\mathbf{c}$. To gain insights, recall the primal problem (11) and its dual (12). Strong duality proved in Theorem 2 under Assumption 2 implies the so-called complementary conditions (Shapiro, 2001), which states that for any optimal primal-dual solution pair $(d\lambda^*, w^*)$, $d\lambda^*(\phi) > 0$ only if $\phi \in \operatorname{argmin}_{\phi \in \Phi} F(\mathbf{w}^*, \phi)$. If we knew $\mathbf{w}^*$ in advance and if we could enumerate all policies in $\operatorname{argmin}_{\phi \in \Phi} F(\mathbf{w}^*, \phi)$, then (11) becomes a finite-dimensional LP and we can then obtain the optimal randomized policy. However, it is not always possible to enumerate all polices in $\operatorname{argmin}_{\phi \in \Phi} F(\mathbf{w}^*, \phi)$. Instead, our approach does not require precise knowledge of $\mathbf{w}^*$, nor the strong duality.

Our approach is to generate a random vector $\tilde{\mathbf{w}}$ and for fixed $\mathbf{w} = \tilde{\mathbf{w}}$ we solve problem (16) to obtain a deterministic policy $\phi_{\mathbf{w}}$. This procedure gives us a randomized policy. The approach is motivated by applying the stochastic gradient descent algorithm (SGD) to solve the dual problem (12), which is reformulated as follows

$$\max_{\mathbf{w} \geq 0} \quad G(\mathbf{w}) := \left\{ \sum_{j \in \mathcal{N}} w_j \cdot \beta_j + \min_{\phi \in \Phi} \left[ \mathrm{E}_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] - \sum_{j \in \mathcal{N}} w_j \cdot \mathrm{E}_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \right] \right\} \quad (25)$$

Specifically, SGD starts from any $\mathbf{w}^{(1)} \geq 0$ and for each $t = 1, 2, \ldots, T$, updates

$$w_j^{(t+1)} = \left[ w_j^{(t)} + \gamma_T \cdot \frac{\partial \hat{G}(\mathbf{w}^{(t)}; \mathbf{D}^{(t)})}{\partial w_j} \right]^+ \quad \forall j \in \mathcal{N}$$

with a step size $\gamma_T$, where $\frac{\partial \hat{G}(\mathbf{w}^{(t)}; \mathbf{D}^{(t)})}{\partial w_j} = \beta_j - R_j \left( s_j \left( \phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)} \right), D_j^{(t)} \right)$ and $\mathbf{D}^{(t)}$ is an independent sample. Then, the expectation of $\frac{\sum_{t=1}^{T} \mathbf{w}^{(t)}}{T}$ will converge to $\mathbf{w}^*$ with an appropriate step size, e.g. $\gamma_T = \frac{1}{\sqrt{T}}$ (Hazan, 2019). Indeed, we can set $\tilde{\mathbf{w}}$ to be the uniform distribution over $\{\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \ldots, \mathbf{w}^{(T)}\}$ to derive our policy.

We present our policy in Algorithm 1. According to step 3 of Algorithm 1, our randomized capacity rationing policy selects a deterministic policy $\phi_{\mathbf{w}^{(t)}}$, $t = 1, 2, 3, \cdots, T$, with probability $1/T$. We refer to this policy as the *Max-Weighted-Service policy*. By Lemma 1, under Assumption 1, for any $\mathbf{w}$, implementing the policy $\phi_{\mathbf{w}}$ in Algorithm 1 only requires solving (16) for the given demand realization $\mathbf{D}$. In the following, for each $t = 1, \ldots, T$, we further denote an i.i.d. copy of $\tilde{\mathbf{D}}$ as $\tilde{\mathbf{D}}^t$, which is also independent of the samples $\{\mathbf{D}^{(1)}, \ldots, \mathbf{D}^{(T)}\}$.

Note that the Max-Weighted-Service policy requires to obtain $T$ samples of the demand distribution. For any policy $\tilde{\phi}$, we denote by $\tilde{\phi}(T)$ if $\tilde{\phi}$ requires $T$ samples of the demand distribution. Then, we call $\tilde{\phi}$ *asymptotically optimal* if and only if

$$\limsup_{T \to \infty} \mathrm{E}_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\tilde{\phi}(T), \mathbf{c}, \tilde{\mathbf{D}}))] \leq \mathrm{Obj} \ (11) \quad \text{and} \quad \liminf_{T \to \infty} \mathrm{E}_{\tilde{\mathbf{D}}} \left[ R_j(s_j(\tilde{\phi}(T), \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j) \right] \geq \beta_j, \quad \forall j \in \mathcal{N}.$$

---

**Algorithm 1** Max-Weighted-Service Policy

---

1: Generate demand samples $\{\mathbf{D}^{(1)}, \mathbf{D}^{(2)}, \ldots, \mathbf{D}^{(T)}\}$ independently from demand distribution $F$, where $T$ is sufficiently large.

2: Starting from $\mathbf{w}^{(1)} = 0$, iteratively generate a random sequence $\{\mathbf{w}^{(2)}, \mathbf{w}^{(3)}, \ldots, \mathbf{w}^{(T+1)}\}$ as follows:

$$w_j^{(t+1)} = \left[ w_j^{(t)} + \gamma_T \cdot \left( \beta_j - R_j \left( s_j \left( \boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)} \right), D_j^{(t)} \right) \right) \right]^+ \tag{26}$$

3: Draw a vector $\tilde{\mathbf{w}}_T$ from $\{\mathbf{w}^{(1)}, \mathbf{w}^{(1)}, \ldots, \mathbf{w}^{(T)}\}$ uniformly at random. Given $\mathbf{w} = \tilde{\mathbf{w}}_T$, adopt the deterministic policy $\boldsymbol{\phi}_{\mathbf{w}}$.

---

where Obj(11) denotes the optimal value of (11). In order to prove the asymptotic optimality of the Max-Weighted-Service policy, it is sufficient to prove that the following two inequalities hold amost surely:

$$\limsup_{T \to \infty} \frac{1}{T} \cdot \sum_{t=1}^{T} \mathrm{E}_{\tilde{\mathbf{D}}^t}[f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t))] \leq \mathrm{Obj}\ (11) \tag{27}$$

and

$$\liminf_{T \to \infty} \frac{1}{T} \cdot \sum_{t=1}^{T} \mathrm{E}_{\tilde{\mathbf{D}}^t} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t) \right] \geq \beta_j, \quad \forall j \in \mathcal{N}. \tag{28}$$

We are now ready to prove the asymptotic optimality of the Max-Weighted-Service policy under the following additional mild assumption.

ASSUMPTION 3. *The support of demand $\tilde{\mathbf{D}}$ is bounded, and there exists a constant $C$ such that for each $j \in \mathcal{N}$, $R_j(s_j, D_j) \leq C$ for each $D_j$ and each $s_j \leq D_j$.*

It is clear that Type I, Type II and Type III service measure functions all satisfy Assumption 3 when the support of demand $\tilde{\mathbf{D}}$ is bounded.

THEOREM 3. *Under Assumption 1 and Assumption 3, if the capacity level $\mathbf{c}$ is feasible and the step size $\gamma_T = T^{-(\frac{1}{2}+\epsilon)}$ for some $\epsilon \in (0, 1/2)$, then the Max-Weighted-Service policy is asymptotic optimal, i.e., (27) and (28) hold almost surely.*

The proof is relegated to Appendix E. Theorem 3 shows the almost surely convergence. If we consider a weaker version of convergence, namely, convergence in expectation, a simple modification of the proof of Theorem 3 also shows the convergence rates in the objective value and the service level constraints. And the results hold under a weaker assumption than Assumption 3.

COROLLARY 2. *Suppose Assumption 1 holds and assume that there exists a constant $C$ such that $E_{\tilde{\mathbf{D}}}[R_j(s_j, \tilde{D}_j)^2] \leq C$ for all $s_j \geq 0$. Then*

$$\frac{1}{T} \cdot \sum_{t=1}^{T} E_{\mathbf{w}(t), \tilde{\mathbf{D}}^t}[f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}(t)}, \mathbf{c}, \tilde{\mathbf{D}}^t))] - Obj\ (11) \leq O(\gamma_T)$$

*and for each $j \in \mathcal{N}$*

$$\beta_j - \frac{1}{T} \cdot \sum_{t=1}^{T} E_{\mathbf{w}(t), \tilde{\mathbf{D}}^t} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}(t)}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t) \right] \leq O(\max\{\sqrt{\frac{1}{T}}, \sqrt{\frac{1}{T\gamma_T}}\})$$

By Corollary 2, we can use different step sizes to prove different convergence rates for the expected allocation cost and the expected service level. For example, by setting $\gamma_T = \frac{1}{\sqrt{T}}$, we get $O(\frac{1}{\sqrt{T}})$ convergence rate for the expected allocation cost and $O(\frac{1}{T^{1/4}})$ convergence rate for expected service levels. By setting $\gamma_T = \frac{1}{T^{1/3}}$, we get $O(\frac{1}{T^{1/3}})$ convergence rate for both the expected allocation cost and expected service levels. If there is no second stage allocation cost, then we can choose an arbitrary $\gamma_T > 0$ and get a convergence rate of $O(\frac{1}{\sqrt{T}})$ on the expected service levels.

### 3.1. Polymatroid and Index Policies

The Max-Weighted-Service policy presented in the previous section is asymptotically optimal for our general model as long as (the very mild) Assumption 1 and Assumption 3 hold. By imposing additional assumptions on the problem structure and the service measure functions, it is possible to obtain additional insights about the policy. The assumption is imposed on a polymatroid structure of the feasible set. The definition of polymatroid is based on submodular set functions (Welsh, 2010). A function $q: 2^{\mathcal{N}} \to R_+$ is called a *submodular* set function if for any $U, V \subseteq \mathcal{N}$, we have

$$q(U) + q(V) \geq q(U \cap V) + q(U \cup V).$$

Moreover, $q$ is non-decreasing if $q(U) \leq q(V)$ for any $U \subseteq V \subseteq \mathcal{N}$. Then a set $Q$ is called a *polymatroid* if there exists a non-decreasing submodular set function $q$ with $q(\emptyset) = 0$ such that

$$Q = \left\{ \mathbf{s} \in \mathbb{R}_+^{\mathcal{N}} \mid \sum_{j \in U} s_j \leq q(U), \quad \forall U \subseteq \mathcal{N} \right\}. \tag{29}$$

The main assumption of this section is the following.

ASSUMPTION 4. *For any given capacity level $\mathbf{c}$ and any realized demand $\mathbf{D}$, the feasible set $Q(\mathbf{c}, \mathbf{D}) = \{\mathbf{s} : \exists \mathbf{y}, (\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})\}$ is a polymatroid.*

A wide range of capacity allocation problems enjoys the polymatorid structure, which includes single-resource pooling, process flexibility, and the generalized W-system ATO problem, as described previously. We relegate the detailed discussions to Appendix F.

In the following, by a straightforward application of a well-known result in polymatroid optimization, we show that an index policy is asymptotically optimal whenever both the allocation cost and the service measure function $R_j(s_j, D_j)$ are linear in $s_j$ for fixed $D_j$, for all $j \in \mathcal{N}$. An *index policy* fulfills the demand of customers according to an index list: the demand of customers with a higher index will be fulfilled as much as possible before fulfilling those with a lower index. Specifically, for each fixed $\mathbf{D}$, the index policy will first lexicographically maximizes $s_j$ according to the priority order while satisfying $\mathbf{s} \in Q(\mathbf{c}, \mathbf{D})$. Once the fulfillment $\mathbf{s}$ is determined, the allocation $\mathbf{y}$ is determined by solving the Max-Weighted-Service problem (16) with $\mathbf{s}$ being fixed. In the single-resource pooling setting, this implies that only one customer is partially satisfied and the ones with higher indices than that customer will be completely satisfied (Alptekinoğlu et al., 2013). However, when multiple resources are involved such that different customers are served by different resources, it is possible that even the customer with the highest index may not be completely fulfilled while a customer with a lower index is fulfilled, since the corresponding resources for the highest index customer may not be enough to fully cover the demand of that customer.

There are two types of index policies, namely responsive and anticipative index policies. An index policy is *responsive* if the index list is constructed *after* demand realization and thus can potentially depend on realized demand, while an index policy is *anticipative* if the index list is constructed *before* demand realization. Both responsive index policies and anticipative index policies can be *deterministic* or *randomized*.

THEOREM 4. *Suppose that the function of the allocation cost can be represented as $\sum_{j \in \mathcal{N}} v_j \cdot s_j$ where $\mathbf{s}$ denotes the fulfillment and for every $j \in \mathcal{N}$, the service measure functions are linear in $s_j$, i.e., $R_j(s_j, D_j) = a_j(D_j) \cdot s_j + b_j(D_j)$, and that Assumption 4 holds. Let $\mathbf{w}$ be the random weight vector generated by Algorithm 1. Denote by $\{i_1, i_1, \ldots, i_n\}$ a permutation of $\{1, 2, \ldots, n\}$ such that $w_{i_j} \cdot a_{i_j}(D_{i_j}) - v_{i_j} \geq w_{i_{j+1}} \cdot a_{i_{j+1}}(D_{i_{j+1}}) - v_{i_{j+1}}$ for all $j = 1, \ldots, n-1$. It is asymptotically optimal to fulfill the demand of the customers in the following way:*

$$s_{i_1}^* = q(\{i_1\}|\mathbf{c}, \mathbf{D})$$
$$s_{i_j}^* = q(\{i_1, i_2, \cdots, i_j\}|\mathbf{c}, \mathbf{D}) - q(\{i_1, i_2, \cdots, i_{j-1}\}|\mathbf{c}, \mathbf{D}), \quad j = 2, \cdots, n$$

The proof is relegated to Appendix F. To conclude this section, we make a couple of remarks.

**Remark 1:** If $R_j(s_j, D_j) = s_j/D_j$ for all $j \in \mathcal{N}$, which corresponds to Type III service constraint, then $a_j(D_j) = 1/D_j$. Therefore, in this case, the customers are fulfilled according to a non-increasing order of $w_j/D_j - v_j$, which depends on the realized demand $\mathbf{D}$. Accordingly, the randomized Max-Weighted-Service policy is a randomized responsive index policy.

**Remark 2:** If $R_j(s_j, D_j) = s_j/\mu_j$ for all $j \in \mathcal{N}$, representing Type II service constraint, then

$a_j(D_j) = 1/\mu_j$. In this case, the customers are fulfilled according to a non-increasing order of $w_j/\mu_j - v_j$, which does not depend on demand realization. Therefore, the randomized Max-Weighted-Service policy is a randomized anticipative index policy. This result unifies and generalizes previous results on capacity allocation with Type II service constraints, including those of Hou et al. (2009) and Zhong et al. (2017) for single-resource pooling, and Lyu et al. (2019) for process flexibility. Specifically, when there is no allocation cost, our policy recovers the policies in the aforementioned papers with the same convergence rate. Moreover, we can show that a randomized anticipative index policy is not just asymptotically optimal, but actually optimal. The detailed discussion is relegated to Appendix F.

## 4. Computing Optimal Capacity Level

Section 3 is focused on characterizing rationing policies for a given capacity level. In this section, we present algorithms to compute optimal capacity levels under Assumption 2. The development of the algorithm relies on the strong duality result in Theorem 2 for fixed $\mathbf{c}$, which then gives rise to a min-max stochastic programming formulation for the original capacity optimization problem (10). To present the formulation, we define for any $\mathbf{w}$ and $\mathbf{c}$,

$$H(\mathbf{w}, \mathbf{c}) = \mathrm{E}_{\tilde{\mathbf{D}}}[h(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})] \tag{30}$$

where for each $\mathbf{D}$,

$$h(\mathbf{w}, \mathbf{c}; \mathbf{D}) = p(\mathbf{c}) + \sum_j w_j \beta_j + g(\mathbf{w}, \mathbf{c}; \mathbf{D}).$$

Recall that $g(\mathbf{w}, \mathbf{c}; \mathbf{D})$ is defined in (16). From the strong duality established in Theorem 2, we have the following result.

THEOREM 5. *Under Assumption 2, problem (10) is equivalent to*

$$\min_{\mathbf{c} \geq 0} \max_{\mathbf{w} \geq 0} \quad H(\mathbf{w}, \mathbf{c}) \tag{31}$$

*in the sense both problems share the same optimal capacity level.*

The proof is relegated to Appendix G. Problem (31) is a minimax stochastic program, for which various optimization algorithms have been developed; see for example Nemirovski et al. (2009). However, in order to guarantee convergence to a globally optimal solution, usually convexity/concavity of the objective function is required.

It is clear that for fixed $\mathbf{c}$ and $\mathbf{D}$, $g(\mathbf{w}, \mathbf{c}; \mathbf{D})$ is concave in $\mathbf{w}$. It then follows immediately that $H(\mathbf{w}, \mathbf{c})$ is concave in $\mathbf{w}$ for any fixed $\mathbf{c}$. However, convexity of $H(\mathbf{w}, \mathbf{c})$ in $\mathbf{c}$ can only be guaranteed with additional assumptions. Lemma 2 below presents one such assumption. The proof can be found in Appendix H.

LEMMA 2. *Under Assumption 2, if the following assumptions hold:*

*(i). The investment cost function $p(\mathbf{c})$ is convex in $\mathbf{c}$,*

*(ii). For all $j \in N$, the service measure function $R_j(s_j, D_j)$ is concave in $s_j$ for any fixed $D_j$,*

*then $H(\mathbf{w}, \mathbf{c})$ is convex in $\mathbf{c}$ for any fixed $\mathbf{w} \geq 0$.*

Obviously, the assumptions of Lemma 2 are satisfied for both Type II and Type III service constraints. Assuming that the optimal capacity level is bounded and $\mathcal{C}$ is a compact convex set containing the optimal capacity level as an interior point, we must have

$$\min_{\mathbf{c} \geq 0} \max_{\mathbf{w} \geq 0} \ H(\mathbf{w}, \mathbf{c}) = \min_{\mathbf{c} \in \mathcal{C}} \max_{\mathbf{w} \geq 0} \ H(\mathbf{w}, \mathbf{c}) = \max_{\mathbf{w} \geq 0} \min_{\mathbf{c} \in \mathcal{C}} \ H(\mathbf{w}, \mathbf{c})$$

When $H(\mathbf{w}, \mathbf{c})$ is convex in $\mathbf{c}$ and concave in $\mathbf{w}$, the maximin problem in the above strong duality relation is the dual problem of our original problem (10) in which $\mathbf{c}$ is a decision variable.

To proceed, we further make the following mild assumption about the objective function.

ASSUMPTION 5. *For any fixed $\varepsilon > 0$, there exists a compact convex set $\mathcal{W}$ in $\mathbb{R}^n$ such that*

$$\max_{\mathbf{w} \geq 0} \min_{\mathbf{c} \in \mathcal{C}} H(\mathbf{w}, \mathbf{c}) - \varepsilon \leq \max_{\mathbf{w} \in \mathcal{W}} \min_{\mathbf{c} \in \mathcal{C}} H(\mathbf{w}, \mathbf{c}) \leq \max_{\mathbf{w} \geq 0} \min_{\mathbf{c} \in \mathcal{C}} H(\mathbf{w}, \mathbf{c}) + \varepsilon$$

It can be easily verified that this assumption in fact holds under Assumption 2 as long as the optimal capacity level is finite. Since our focus in this section is to apply an existing algorithm to solve our problem, we skip the verification of this assumption.

We now apply the mirror descent stochastic approximation (SA) algorithm of Juditsky and Nemirovski (2011) to solve the minimax stochastic program (31). The presentation of the algorithm requires the following notations. Let $\|\cdot\|$ be a general norm defined in $\mathbb{R}^{2n}$, with $\|x\|_* = \sup_{\|v\| \leq 1} v^T x$ being its dual norm. Let $l : X = \mathcal{C} \times \mathcal{W} \to \mathbb{R}$ be a distance-generating function. If $l(\cdot)$ is convex and continuous on $X$, the set

$$X^0 = \{x \in X : \ \exists u \in \mathbb{R}^{2n} \ \text{s.t.} \ x \in \mathrm{argmin}_{v \in X}[u^T v + l(v)]\}$$

is convex. Suppose $l(\cdot)$ is continuously differentiable and strongly convex on $X^0$ with parameter 1 with respect to $\|\cdot\|$, i.e.,

$$(x' - x)^T (\nabla l(x') - \nabla l(x)) \geq \|x' - x\|^2, \quad \forall x', x \in X^0.$$

The prox-function is defined by

$$V(x, z) = l(z) - [l(x) + \nabla l(x)^T (z - x)]$$

and the prox mapping is defined by $\mathcal{P}_x : \mathbb{R}^{2n} \to X^0$ such that

$$\mathcal{P}_x(u) = \mathrm{argmin}_{z \in X}\{u^T(z - x) + V(x, z)\}.$$

There are many ways to choose distance generating functions: for example $l(x) = \sum_{k=1}^{2n} x_k \log(x_k)$ or $l(x) = \frac{1}{2}\|x\|_2^2$. In the following analysis, we adopt the Euclidean norm $\|\cdot\|_2$. For notational brevity, we define the tuple $\theta = [\mathbf{w}^T, \mathbf{c}^T]^T$ and denote $\partial h(\theta; \mathbf{D}) = \partial h(\mathbf{w}, \mathbf{c}; \mathbf{D})$, which is an unbiased estimator of $\partial H(\mathbf{w}, \mathbf{c})$ represented by

$$\partial h(\mathbf{w}, \mathbf{c}; \mathbf{D}) = \begin{bmatrix} \partial_{\mathbf{c}} h(\mathbf{w}, \mathbf{c}; \mathbf{D}) \\ -\partial_{\mathbf{w}} h(\mathbf{w}, \mathbf{c}; \mathbf{D}) \end{bmatrix}.$$

Then the mirror descent SA algorithm is presented in Algorithm 2.

---

**Algorithm 2** SA Algorithm Juditsky and Nemirovski (2011)

---

**Input:** initial point $\theta^{(1)}$, time horizon $T$, positive step size $\{\gamma^{(t)}\}_{t=1}^T$, and a sequence $\{\mathbf{D}^{(t)}\}_{t=1}^T$, which is a sequence of samples of $\tilde{\mathbf{D}}$.

**Output:** sequence $\{\theta^{(t)}\}_{t=1}^T$.

    **for** $t = 1, ..., T$ **do**

        $\theta^{(t+1)} = \mathcal{P}_{\theta^{(t)}}(\gamma^{(t)} \cdot \partial h(\theta^{(t)}, \mathbf{D}^{(t)}))$

    **end for**

---

Proposition 1.7 of Juditsky and Nemirovski (2011) implies that Algorithm 2 converges to an optimal capacity level when $E[\|\partial h(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})\|_2^2]$ is bounded and the step size $\{\gamma^{(t)}\}_{t=1}^T$ satisfies

$$\sum_{t=1}^T \gamma^{(t)} \to \infty \quad \text{and} \quad \frac{\sum_{t=1}^T (\gamma^{(t)})^2}{\sum_{t=1}^T \gamma^{(t)}} \to 0 \quad \text{as} \quad T \to \infty.$$

A common choice of step size is the constant step size $\gamma^{(t)} = \frac{\delta}{\sqrt{T}}$ for $t = 1, 2, \ldots, T$, where $\delta > 0$ is a parameter. One could also choose the step size $\gamma^{(t)} = \frac{\delta}{\sqrt{t}}$ for $t = 1, 2, \ldots, T$, which does not require a fixed total number of iterations in advance.

The SA algorithm can be directly applied to solve problem (31) whenever $H(\mathbf{w}, \mathbf{c})$ is convex in $\mathbf{c}$ (it is always concave in $\mathbf{w}$). We report its numerical performance in the next section when it is applied to problems with Type II and Type III service constraints.

However, when the service constraints are of Type I, $H(\mathbf{w}, \mathbf{c})$ is not guaranteed to be convex in $\mathbf{c}$ in general and $\partial h(\mathbf{w}, \mathbf{c}; \mathbf{D})$ may not be well defined. In this case, we propose a heuristic to solve problem (31) approximately. The details are provided below. Note that the (super)-gradient of $h(\mathbf{w}, \mathbf{c}; \mathbf{D})$ over $\mathbf{w}$ is given by

$$\partial_w h(\mathbf{w}, \mathbf{c}; \mathbf{D}) = \beta - \mathbf{z}^* := (\beta_j - z_j^*, \ j \in N) \tag{32}$$

where $\mathbf{z}^*$ denotes an optimal solution to the problem below, which is a reformulation of the Max-Weighted-Service problem for Type I service constraints,

$$g(\mathbf{w}, \mathbf{c}, \mathbf{D}) = \min_{\mathbf{y}, \mathbf{s}, \mathbf{z}} f(\mathbf{y}) - \sum_{j \in \mathcal{N}} w_j z_j \ \text{ s.t. } (\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D}), \ s_j \geq z_j D_j, \ z_j \in \{0, 1\} \ \forall j \in \mathcal{N}. \tag{33}$$

To get an approximation of $\partial_c h(\mathbf{w}, \mathbf{c}; \mathbf{D})$, we relax the integer constraint in (33) and get an LP relaxation. Let $\varphi$ be the coefficient of $\mathbf{c}$ in the objective value of the dual of the LP relaxation, and $\varphi^*$ be the value of $\varphi$ in an optimal dual solution. Define

$$\partial_c \tilde{h}(\mathbf{w}, \mathbf{c}; \mathbf{D}) = \nabla p(\mathbf{c}) + \varphi^*$$

as an approximation of $\partial_c h(\mathbf{w}, \mathbf{c}; \mathbf{D})$. Thus,

$$\partial \tilde{h}(\mathbf{w}, \mathbf{c}; \mathbf{D}) = \begin{bmatrix} \partial_c \tilde{h}(\mathbf{w}, \mathbf{c}; \mathbf{D}) \\ -\partial_w h(\mathbf{w}, \mathbf{c}; \mathbf{D}) \end{bmatrix} = \begin{bmatrix} \nabla p(\mathbf{c}) + \varphi^* \\ -\beta + \mathbf{z}^* \end{bmatrix}$$

will be used to approximate $\partial h(\mathbf{w}, \mathbf{c}; \mathbf{D})$ in Algorithm 2. We refer to this heuristic as the SA heuristic. However, it is not guaranteed that such a heuristic can always obtain a feasible capacity level. In practice, if the solution of this heuristic is infeasible, we can apply a problem-specific procedure to increase the capacity level. The details and the numerical performance of this method will be reported in the next section.

## 5. Numerical Results

We now present numerical results to demonstrate the performance of the SA algorithm and the SA heuristic proposed in the previous section. In our numerical study, We choose the distance-generating function $l(x) = \frac{1}{2}\|x\|_2^2$. We set the transportation cost to 0 in the numerical experiments. The SA algorithm is guaranteed to converge to optimal capacity level if the service constraints are convex, e.g., they are of Type II or Type III. In Section 5.1, we numerically illustrate the rate of convergence for the SA algorithm when the service levels are of Type II or Type III. From Section 5.2 to Section 5.4, we investigate the effectiveness of our algorithm under Type I service level in different applications. For Type I service constraints, We focus on the accuracy of the SA Heuristic. We study inventory pooling problems in Section 5.2, flexible production problems in Section 5.3 and the assemble-to-order problems in Section 5.4.

### 5.1. Convergence Rate of SA Algorithm under Type II and Type III Service Levels

In this section, we apply the SA algorithm to compute the optimal capacity level under Type II and Type III service levels in different applications, including inventory pooling problems, flexible production problems, and assemble-to-order problems. We aim to evaluate the rate of convergence of the SA algorithm under different applications, system sizes, and service levels. For the inventory pooling problem, we run the SA algorithm in the $1 \times 10$, $1 \times 20$, and $1 \times 50$ systems, where the demand follows a truncated normal distribution with a mean of 1 and a standard deviation of 0.3. For the flexible production problem, we evaluate the SA algorithm in the $10 \times 10$, $20 \times 20$, and $50 \times 50$ long-chain system (illustrated in Figure 1), where the demand follows the same distribution

as in the inventory pooling problem. In the Assemble-To-Order problem, we test the algorithm under the $11 \times 10$, $21 \times 20$, and $51 \times 50$ generalized W-systems (illustrated in Figure 1) with the same demand distribution. The target service levels are the same among all products, and all the applications are replicated for service levels being 0.85, 0.90 and 0.95. The structures of different systems are presented in Figure 1.

**Figure 1** Structure of Inventory Pooling, Long Chain, and the ATO System in Section 5.1
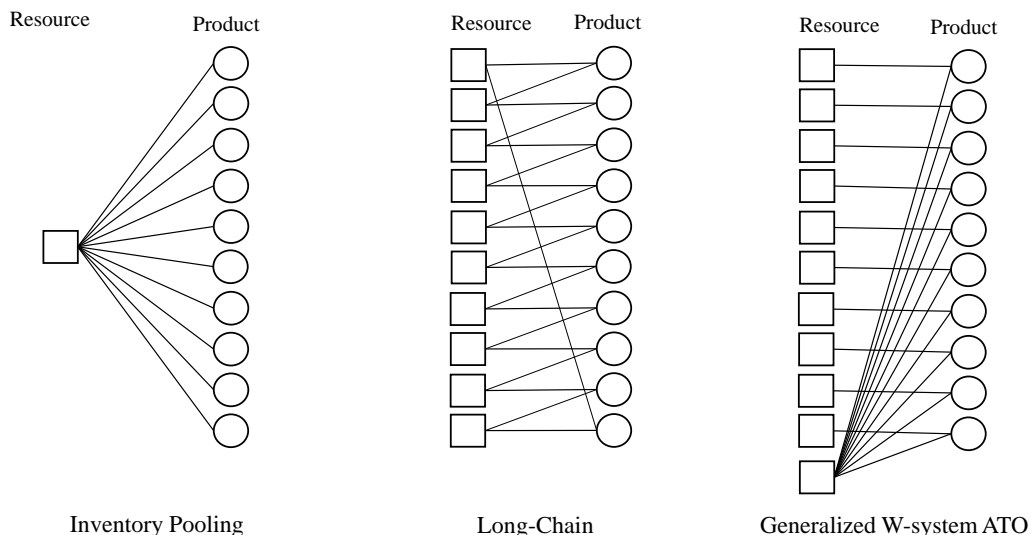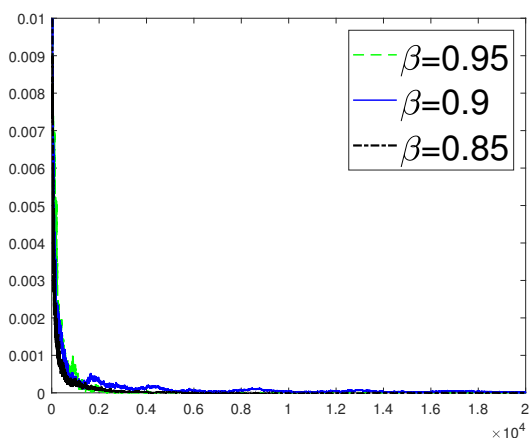


Inventory Pooling    Long-Chain    Generalized W-system ATO

**Table 1** The Convergence Time for SA Algorithm under Type II and Type III Service Levels

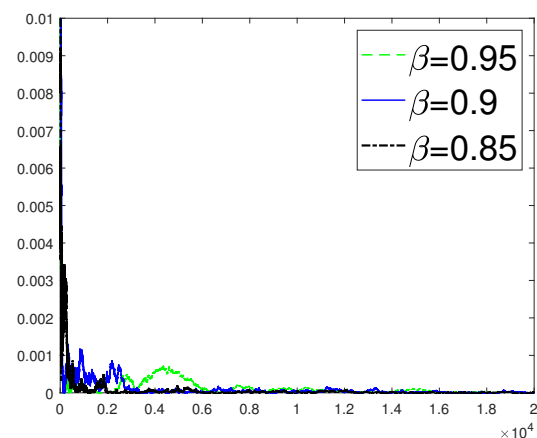| Type II | Inventory Pooling | | | Long Chain | | | Assemble-To-Order | | |
|---|---|---|---|---|---|---|---|---|---|
| n  SL | 0.85 | 0.9 | 0.95 | 0.85 | 0.9 | 0.95 | 0.85 | 0.9 | 0.95 |
| 10 | 1010 | 3107 | 4936 | 27496 | 83824 | 42866 | 50151 | 61770 | 84655 |
| 20 | 588 | 1620 | 2531 | 35899 | 93035 | 53668 | 67292 | 75103 | 90597 |
| 50 | 461 | 981 | 2219 | 26182 | 93993 | 55964 | 70236 | 72824 | 81097 |
| Type III | Inventory Pooling | | | Long Chain | | | Assemble-To-Order | | |
| n  SL | 0.85 | 0.9 | 0.95 | 0.85 | 0.9 | 0.95 | 0.85 | 0.9 | 0.95 |
| 10 | 4806 | 8618 | 7572 | 24705 | 30401 | 43187 | 29708 | 31596 | 71979 |
| 20 | 15760 | 14493 | 6280 | 39229 | 40189 | 17444 | 36352 | 44435 | 81039 |
| 50 | 35444 | 32342 | 26459 | 10802 | 31518 | 22951 | 24554 | 38556 | 55206 |

For each instance, we run the SA algorithm for $T = 200000$ periods. We observe that the SA algorithms converge within 100000 periods. To better approximate the optimal objective value, we approximate it by the average value of the objectives in the first 200000 periods, i.e. $\mathrm{obj}^* = \frac{1}{T}\sum_{\tau=1}^{T}\mathrm{obj}(\tau)$ with $T = 200000$. To evaluate the convergence rate in each instance, we find the largest $t$ such that $|\frac{1}{t}\sum_{\tau=1}^{t}\mathrm{obj}(\tau) - \mathrm{obj}^*| \geq 0.1\% \cdot \mathrm{obj}^*$. Then we refer to $t$ as the convergence time. In Table 1, we present the convergence time in different applications, system sizes and service

levels. Table 1 indicates that under the Type II service level, the SA algorithm converges the fastest for the inventory pooling problems. Moreover, increasing the system size will not always increase the convergence time but sometimes reduce the convergence time in inventory pooling problems. Notice that here the convergence time is the number of iterations or periods the SA algorithm takes before convergence. For the actual running time in seconds, it will increase in the system size. We note that increasing the service level often increases the convergence time. Besides, for inventory pooling problems, the SA algorithm takes more periods to converge under the Type III service level than that under the Type II service level. For the Long Chain system and the ATO systems, the convergence time is comparable for the Type II and Type III service levels. In most instances, the convergence time is within 50000 periods, and only a few instances take more periods to converge but still less than 100000 periods. For the optimal capacity we

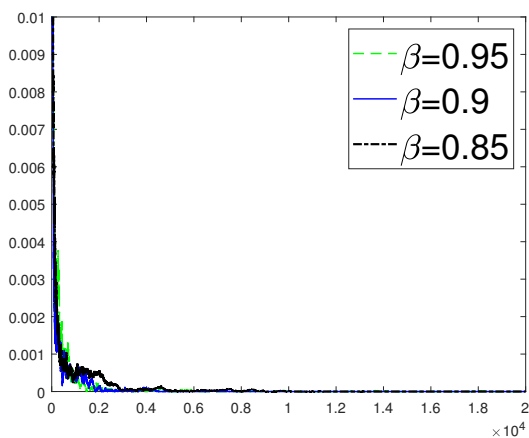**Figure 2**     The Gap between the Target Service Level and the Achieved Service Level under Corollary 1



(a) Long Chain System under Type II Service Level

(b) Long Chain System under Type III Service Level

(c) ATO System under Type II Service Level

(d) ATO System under Type III Service Level

solved by the SA heuristic, we adopt Corollary 1 to check the feasibility of the capacity level. We

plot the sum of all the positive part of the difference between the target service level and the achieved service level under Corollary 1 among all products, divided by the number of products, i.e., $\frac{1}{n}\sum_{j\in\mathcal{N}}\left(\beta_j - \frac{1}{t}\sum_{\tau=1}^{t} R_j\left(s_j\left(\boldsymbol{\phi}_{\mathbf{w}(\tau)}, \mathbf{c}, \mathbf{D}^{(\tau)}\right), D_j^{(t)}\right)\right)^+$ for the flexible production problems and the assemble-to-order problems with 20 products in Figure 2. The diminishing of the sum of the positive part of the gap implies that the target service levels are achieved for all products within 10000 periods.

We have illustrated the convergence rate of the SA algorithm in computing the optimal capacity level, and the convergence rate of the service level by Corollary 1. Once the capacity level

**Figure 3**     The Gap between the Target Service Level and the Achieved Service Level under Algorithm 1
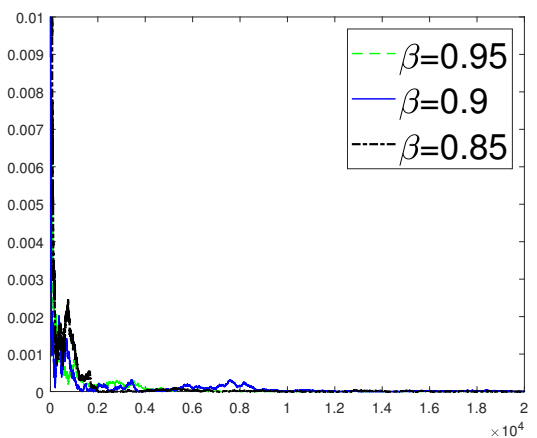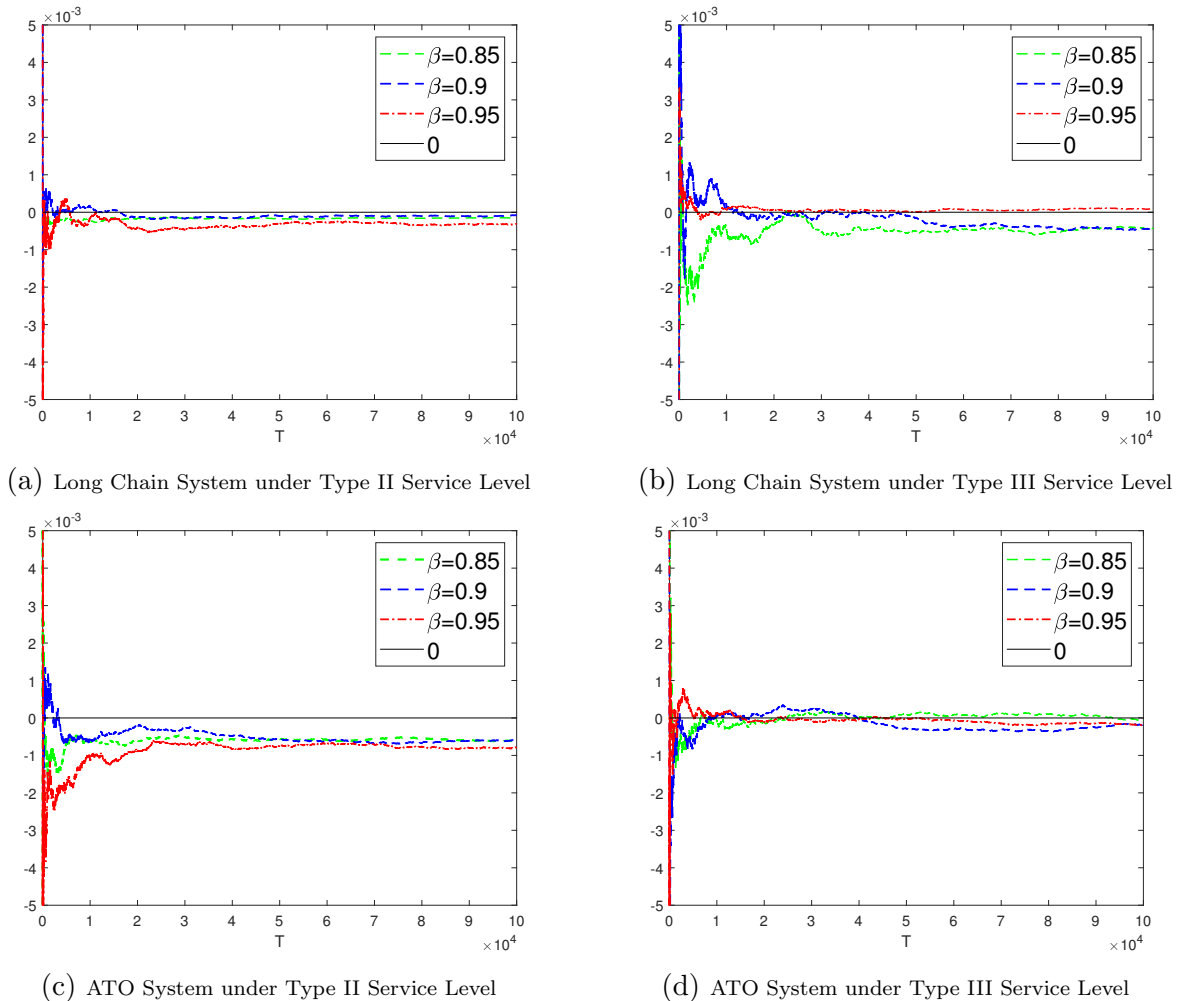


(a) Long Chain System under Type II Service Level

(b) Long Chain System under Type III Service Level

(c) ATO System under Type II Service Level

(d) ATO System under Type III Service Level

is fixed, we can use Algorithm 1 (the Max-Weighted-Service policy) to allocate the resources to achieve the target service level. We now evaluate the convergence of Algorithm 1 in service levels for the solved capacity by the SA algorithm under different settings. We draw $T = 100000$

number of independent demand samples, and then generate the $\mathbf{w}$ sequence according to Algorithm 1. After that, we draw a random vector $\mathbf{w}$ from the generated sequence, and adopt the allocation policy $\phi_{\mathbf{w}}$ to a new independently generated demand sample. We evaluate the gap in service level as the number of new demand samples increases. Specifically, Figure 3 depicts $\frac{1}{n} \sum_{j \in \mathcal{N}} \left( \beta_j - \frac{1}{t} \sum_{\tau=1}^{t} R_j \left( s_j \left( \phi_{\mathbf{w}^{(\tau)}}, \mathbf{c}, \mathbf{D}^{(\tau)} \right), D_j^{(t)} \right) \right)$ for the flexible production problems and the assemble-to-order problems with 20 products. It is illustrated in Figure 3 that the service level achieves the target service level at a fast rate.

## 5.2. Inventory Pooling with Type I Service Constraints

From now on, we will focus on calculating the capacity level for Type I service level. In this subsection, we consider an inventory pooling example with Type I service constraints. There are $n = 10$ customers and the demands are i.i.d. normal distributions with a mean of 10 units. The standard deviation is set to be either 3 or 5 units. For each value of standard deviation, we test the performance of our heuristic with six sets of target service levels presented in Table 2. For example, in Exp1, the service levels vary between 71% and 89% with an average of 80%. In Exp2, all service levels are equal to 80%.

**Table 2**    Six Sets of Service Levels for All Customers

|      | $\beta_1$ | $\beta_2$ | $\beta_3$ | $\beta_4$ | $\beta_5$ | $\beta_6$ | $\beta_7$ | $\beta_8$ | $\beta_9$ | $\beta_{10}$ |
|------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Exp1 | 0.71  | 0.73  | 0.75  | 0.77  | 0.79  | 0.81  | 0.83  | 0.85  | 0.87  | 0.89  |
| Exp2 | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   |
| Exp3 | 0.85  | 0.85  | 0.85  | 0.85  | 0.85  | 0.85  | 0.85  | 0.85  | 0.85  | 0.85  |
| Exp4 | 0.855 | 0.865 | 0.875 | 0.885 | 0.895 | 0.905 | 0.915 | 0.925 | 0.935 | 0.945 |
| Exp5 | 0.9   | 0.9   | 0.9   | 0.9   | 0.9   | 0.9   | 0.9   | 0.9   | 0.9   | 0.9   |
| Exp6 | 0.95  | 0.95  | 0.95  | 0.95  | 0.95  | 0.95  | 0.95  | 0.95  | 0.95  | 0.95  |

For each set of target service levels, we apply the SA heuristic to compute the total inventory. After adopting SA heuristic, we test the feasibility of the current capacity level by Corollary 1. If the capacity is infeasible, we increase the capacity to a feasible level using a bisection method. As a benchmark, we also apply the greedy algorithm of Alptekinoğlu et al. (2013) together with a bisection search to Exp2, Exp3, Exp 5, Exp6; this algorithm is theoretically optimal when the demands are i.i.d. and service levels are equal for all customers.

As is reported in Table 3, the SA heuristic performs extremely well; the computed inventory is always within 1% of the optimal solution for all cases. The optimal solutions for Exp 1 and Exp 4 are not included in Table 3, since the greedy algorithm of Alptekinoğlu et al. (2013) does not directly apply when service levels are not uniform. For these two sets of parameters, we evaluate the performance of the SA heuristic with the help of the following proposition.

**Table 3**   Capacity Level of Inventory Pooling: Optimal Algorithm vs. SA Heuristic

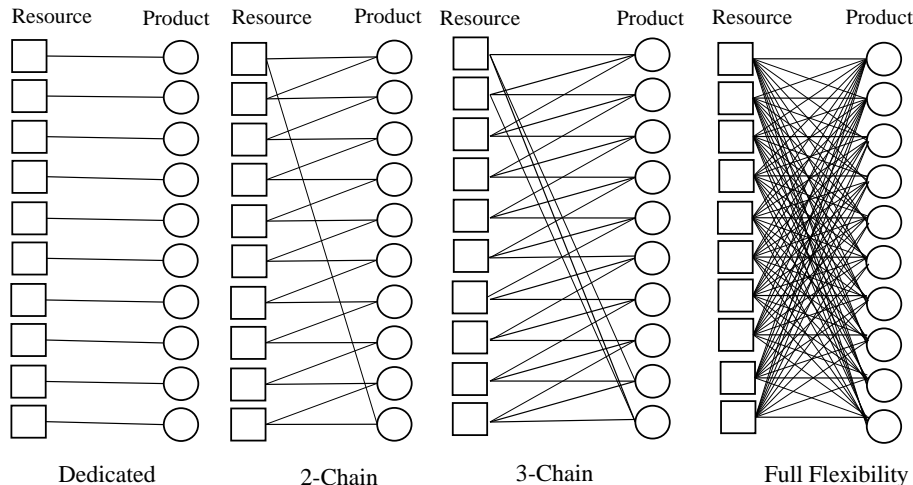|  | N(10,3) | | N(10,5) | |
|---|---|---|---|---|
|  | Optimal | SA Heuristic | Optimal | SA Heuristic |
| Exp1 | - | 78.7213 | - | 75.6031 |
| Exp2 | 78.5471 | 78.5612 | 74.7758 | 75.4424 |
| Exp3 | 85.2919 | 85.3004 | 82.7857 | 83.4993 |
| Exp4 | - | 92.6152 | - | 92.4382 |
| Exp5 | 92.5104 | 92.5202 | 91.8715 | 92.3303 |
| Exp6 | 100.9395 | 100.9757 | 103.4066 | 103.8292 |

PROPOSITION 1. *Consider inventory pooling with $n$ customers with i.i.d. demand distributions. The optimal capacity level with differentiated service levels $\beta = (\beta_1, \beta_2, \ldots, \beta_n)$ is greater than or equal to that with the uniform target service level $\hat{\beta} = \frac{1}{n} \sum_{j \in \mathcal{N}} \beta_j$ for all customer $j \in \mathcal{N}$.*

The proof is relegated to Appendix J. Therefore, the optimal inventory level required in Exp1 should be at least as high as that in Exp2. However, Table 3 shows that the computed inventory levels in these two cases are very close to each other, and both of them should be close to their respective optimal solutions. Similar conclusions can be made for Exp4.

### 5.3. Flexible Production with Type I Service Constraints

In this subsection, we consider a flexible production problem with $m = n = 10$ under Type I service level. The setup for the 10 products/customers is the same as that in the previous subsection. We will consider different flexibility designs including the dedicated design, long chain, 3-chain, and full flexibility design. The designs are illustrated in Figure 4. A formal definition of these designs can be found in Chou et al. (2010); Wang and Zhang (2015). After adopting the SA heuristic, we check the feasibility of the current capacity level, denoted by **c**, by Corollary 1. If infeasible, we apply the Max-Weighted-Service policy with a fixed $T$ to determine which resource capacity should be increased. Denote $j(\mathbf{c})$ as the product whose service level achieved by the Max-Weighted-Service policy is the farthest below the target. We use a bisection procedure to increase the capacity of each resource serving $j(\mathbf{c})$ by the same amount so that $j(\mathbf{c})$ can achieve its target service level. We repeat the above procedures to increase the capacity of the resources until the capacity level is feasible. The numerical results are presented in Table 4 and Table 5.

Notice that under full flexibility, the problem reduces to the single-resource pooling case. Therefore, its optimal capacity level can be computed as reported in the previous subsection. With the dedicated design, the optimal capacity level can be computed either via a bisection approach or from the cumulative distribution of the demand. Thus, the capacity levels reported in Table 4 and Table 5 are optimal for the dedicated design and the full flexibility design. The capacity levels for the long chain and 3-chain designs are computed by the SA heuristic.

**Figure 4**    Flexibility structures



**Table 4**    Total capacity of different flexibility structures with 10 customers, i.i.d. demands $N(10, 3)$

| Type I | Dedicated | Long Chain | 3-Chain | Full Flexibility |
|---|---|---|---|---|
| Exp1 | 125.8221 | 82.4799 | 78.8428 | 78.6397 |
| Exp2 | 125.2486 | 82.3647 | 78.8018 | 78.5612 |
| Exp3 | 131.0930 | 87.9725 | 85.5809 | 85.3004 |
| Exp4 | 139.0002 | 94.3152 | 92.8290 | 92.5664 |
| Exp5 | 138.4465 | 94.3217 | 92.7564 | 92.5202 |
| Exp6 | 149.3456 | 102.4108 | 101.1873 | 100.9757 |

**Table 5**    Total capacity of different flexibility structures with 10 customers, i.i.d. demands $N(10, 5)$

| Type I | Dedicated | Long Chain | 3-Chain | Full Flexibility |
|---|---|---|---|---|
| Exp1 | 143.0368 | 85.5178 | 75.8044 | 75.6022 |
| Exp2 | 142.0811 | 85.2705 | 75.6913 | 75.4424 |
| Exp3 | 151.8217 | 92.2134 | 83.7568 | 83.4994 |
| Exp4 | 165.0003 | 100.6091 | 92.7365 | 92.4172 |
| Exp5 | 164.0776 | 100.4087 | 92.5862 | 92.3303 |
| Exp6 | 182.2427 | 112.2169 | 104.1527 | 103.8292 |

We observe that the performance of 3-chain is almost as good as the full flexibility design for Type I service constraints. This phenomenon is well-known when total capacity is fixed and the objective is to optimize demand fulfillment (Chou et al., 2011; Chen et al., 2015; Simchi-Levi and Wei, 2015, 2012; Désir et al., 2016), and is also observed when the service constraints are of Type II (Lyu et al., 2019). Table 4 shows that under Type I constraints, the long chain design can achieve most of the pooling benefit of full flexibility, and the improvement from 3-chain to full flexibility is negligible. This also provides a verification of the numerical performance of our SA heuristic for the long-chain and 3-chain designs with Type I service constraints.

### 5.4. Heuristic for Assemble to Order

In Section 5.2 and Section 5.3, we have demonstrated that the SA heuristic performs well under Type I service level when applied to inventory pooling and process flexibility. However, for ATO problems, the capacity computed by SA Heuristic is at least 30% higher than the optimal level, as reported in Table 8. This motivates us to develop another local search heuristic that is based on the Max-Weighted-Service policy for fixed capacity levels.

Our local search heuristic requires a reasonably good lower bound for the optimal capacity to serve as a starting point, which can be obtained as follows. For each $i$, let $U_i$ be the set of products that share component $i$; if $i$ is product-specific, then $U_i$ contains only one product. This decomposes the problem into a set of subproblems, each of which is a single-resource allocation problem. We can either use the SA heuristic or bisection search to compute the optimal capacity level $c_i^L$ of component $i$ to achieve the target service levels of the products in $U_i$. Then it is clear that $\mathbf{c}^L = (c_1^L, \cdots, c_m^L)$ must be a lower bound for the original problem.

Our heuristic then increases the capacity level to achieve feasibility. If the current capacity level $\mathbf{c}$ is infeasible, then we apply the Max-Weighted-Service policy with a fixed $T$ to determine the achieved service level, denoted by $\beta_j(\mathbf{c})$, for each $j \in \mathcal{N}$. Denote by $j(\mathbf{c})$ the product that has the largest gap between the target and the achieved service levels $\beta_j - \beta_j(\mathbf{c})$ among all $j \in \mathcal{N}$. Denote $S_{j(\mathbf{c})}$ as the set containing all the components that are required by product $j(\mathbf{c})$. We then carefully choose a subset of components in $S_{j(\mathbf{c})}$ to increase their capacities. Specifically, re-index all components in $S_{j(\mathbf{c})}$ such that $\delta_1 \geq \delta_2 \geq \cdots \geq \delta_{|S_{j(\mathbf{c})}|}$ where $\delta_i$ is the number of products in $\mathcal{N}$ connected to component $i$. Then let $Z_{j(\mathbf{c})}$ be a set of subsets of $S_{j(\mathbf{c})}$ such that any $\hat{I} \in Z_{j(\mathbf{c})}$ is of the form $\{i, i+1, \cdots, i+r\}$ for some $i \in S_{j(\mathbf{c})}$ and some integer $r \geq 0$ where $i+k$ denotes $i+k$ module $|S_{j(\mathbf{c})}|$ if $i+k > |S_{j(\mathbf{c})}|$. The total number of subsets in $Z_{j(\mathbf{c})}$ is $|S_{j(\mathbf{c})}|^2$. For each $\hat{I} \in Z_{j(\mathbf{c})}$ we define $M(\hat{I})$ as the ratio of the marginal decrease of the total gap of the service levels after increasing the capacity of each component in $\hat{I}$ by one unit, over $|\hat{I}|$. That is,

$$M(\hat{I}) := \frac{\sum_{j \in \mathcal{N}} (\beta_j - \beta_j(\mathbf{c}))^+ - \sum_{j \in \mathcal{N}} \left(\beta_j - \beta_j(\mathbf{c} + \sum_{i \in \hat{I}} e_i)\right)^+}{|\hat{I}|} \tag{34}$$

where $(a)^+ = \max\{a, 0\}$, $e_i$ is an $m$-dimensional vector with 1 as the $i$-th component and 0 as other components, and the service level $\beta_j(\mathbf{c} + \sum_{i \in \hat{I}} e_i)$ for each $j \in \mathcal{N}$ is again computed by applying the Max-Weighted-Service policy with the same $T$ under the capacity level $\mathbf{c} + \sum_{i \in \hat{I}} e_i$. Denote $\hat{I}^* = \operatorname{argmax}_{\hat{I} \in Z_{j(\mathbf{c})}} M(\hat{I})$, we then increase the capacity of each component in $\hat{I}^*$ by one unit to obtain a new capacity level. We repeat the above procedures until the obtained capacity level is feasible. Our heuristic is formally presented in Algorithm 3.

For the numerical experiments, we consider the following two types of configurations of the ATO system: the generalized W-systems and non-W-systems. We use $W(n)$ and $NW(n)$ to denote the

---

**Algorithm 3** Local Search Heuristics

---

1: Compute a lower bound $\mathbf{c}^L$ and denote $\mathbf{c} = \mathbf{c}^L$.

2: Apply the Max-Weighted-Service policy in Algorithm 1 with a fixed $T$ to compute the achieved service level under $\mathbf{c}$, denoted by $\beta_j(\mathbf{c})$, for each $j \in N$.

3: **while** the capacity level $\mathbf{c}$ is infeasible, **do**

4:    Denote $j(\mathbf{c}) = \mathrm{argmax}_{j \in \mathcal{N}}(\beta_j - \beta_j(\mathbf{c}))$.

5:    Denote $S_{j(\mathbf{c})}$ as a set of all components required by product $j(\mathbf{c})$ and re-index all components in $S_{j(\mathbf{c})}$ in non-increasing order of the number of products connected to the component.

6:    Denote $Z_{j(\mathbf{c})}$ as the set of subsets of $S_{j(\mathbf{c})}$ such that any $\hat{I} \in Z_{j(\mathbf{c})}$ is of the form $\{i, i + 1, \cdots, i + r\}$ for some $i \in S_{j(\mathbf{c})}$ and some integer $r \geq 0$ where $i + k$ denotes $i + k$ module $|S_{j(\mathbf{c})}|$ if $i + k > |S_{j(\mathbf{c})}|$..

7:    Compute $\hat{I}^* = \mathrm{argmax}_{\hat{I} \in Z_{j(\mathbf{c})}} M(\hat{I})$ where $M(\hat{I})$ is defined in (34).

8:    Define $\hat{\mathbf{c}} = \mathbf{c} + \sum_{i \in \hat{I}^*} e_i$ and denote $\mathbf{c} = \hat{\mathbf{c}}$.

9: **end while**

10: Output: $\mathbf{c}$.

---

generalized W-system and non-W-system, respectively, with $n$ products. In NW($n$), there are $2n$ components and for each $i = 1, 2, \ldots, n$, one unit of product $i$ requires one unit of component $i$, component $n + i$ and component $n + i + 1$ where component $2n + 1$ denotes component $n + 1$. Moreover, in both W($n$) and NW($n$), the Type I service level of each product is set to be 95%. In each configuration, the demand of each product is independent normal distributions $N(\mu, \sigma)$ and is rounded to non-negative integers, where $\mu$ (resp. $\sigma$) is sampled from a uniform distribution over $\{9, 10, 11\}$ (resp. $\{2, 3, 4\}$). We implement Algorithm 3 on problem instances W($n$) and NW($n$) for $n = 5, 10, 15, \ldots, 50$.

In order to evaluate the performance of Algorithm 3, we further use the sample average approximation (SAA) method as a benchmark. In SAA, we generate $K$ demand scenarios, denoted as $\{\mathbf{D}^{(1)}, \ldots, \mathbf{D}^{(K)}\}$, and for each $k = 1, \ldots, K$, we introduce a binary variable $z_j^{(k)}$ for each $j \in \mathcal{N}$ to indicate whether the demand of product $j$ is fully satisfied under scenario $\mathbf{D}^{(k)}$. Then, we solve the following integer programming:

$$\min_{\mathbf{c}, \mathbf{z}} \sum_{i \in \mathcal{M}} c_i \tag{35}$$

$$\text{s.t.} \sum_{j \in \mathcal{N}} A_{ij} \cdot D_j^{(k)} \cdot z_j^{(k)} \leq c_i, \quad \forall\, i \in \mathcal{M},\, \forall\, k = 1, \ldots, K$$

$$\sum_{k=1}^{K} z_j^{(k)} \geq \beta_j \cdot K \quad \forall\, j \in \mathcal{N}$$

$$c_i \geq 0, \quad \forall\, i \in \mathcal{M},\; z_j^{(k)} \in \{0,1\} \quad \forall\, j \in \mathcal{N}, \forall\, k = 1, \cdots, K$$

Notice that the number of binary variables in (35) is $n \cdot K$. In our experiments, the computer we use was not able to solve this integer program when $n \geq 35$ and $K = 500$.

However, when $K$ is small (no more than 500), the optimal capacity obtained by solving (35) may not be feasible under the original demand distribution. Therefore, in order for SAA to be a useful benchmark, we first consider replacing the original distribution of each problem instance with a uniform distribution over the $K$ samples $\{\mathbf{D}^{(1)}, \ldots, \mathbf{D}^{(K)}\}$ for $K = 50, 100, 500$. Then, the solution obtained by SAA must be feasible under this distribution.

We remark that once the samples are fixed, the indicator $\mathbf{z}^{(k)}$ determined by (35) is deterministic for each scenario $\mathbf{D}^{(k)}$. Thus, even when the true distribution is the uniform distribution over $\{\mathbf{D}^{(1)}, \ldots, \mathbf{D}^{(K)}\}$, the optimal solution of (35) only characterizes the optimal deterministic policy and the corresponding optimal objective value can be higher than the optimal total capacity under the optimal randomized policy, especially for small $K$. Therefore, in some of the instances, SAA produces a higher capacity than our Algorithm 3.

For the reasons discussed above, we implement Algorithm 3 for each problem instance assuming the true distribution is a uniform distribution over the sampled demand $\{\mathbf{D}^{(1)}, \ldots, \mathbf{D}^{(K)}\}$. The numerical results are reported in Table 6 and Table 7, for the generalized W-system and non-W-system, respectively. Note that for all instances, the gap between the total capacity obtained by SAA and the total capacity obtained by Algorithm 3 is within 2.51%. For each instance, we present the number of times that we increase the capacity in Algorithm 3 in the 'step' column. In the W-system, the number of steps for Algorithm 3 is upper bounded by 25 and the run time is below 1.5 hours. For the non-W-system, since each product requires 3 components, it will require more computation efforts for Algorithm 3 to reach feasibility. However, even when $n = 35, 40, 45, 50$ and $K = 500$, Algorithm 3 still finds a feasible capacity level within 6 hours, while the integer programming (35) for SAA is out of the memory (OOM) of 16 GB[1], and thus no result is obtained.

Recall that results presented in Table 6 and Table 7 treat the sampled distribution as the true distribution. Therefore, the computed capacity may not be feasible for the problem under the original (normal) distribution. In fact, we find that none of the SAA solutions is feasible under the original distribution. For example, consider the W-system with $n = 25$ and equal target service level 95%. When the sample size is 50, with the capacity computed by SAA, there are 9 products that can not achieve a 95% service level. Two of them can only achieve a service level between 89% and 91%. When the sample size is 500, five products can not achieve the target several levels, but are very close (higher than 94.3%).

---

[1] The CPUs we use are 2x Intel Xeon Platinum 8268 24C 205W 2.9GHz Processor, and we assign 16 GB memory for each program.

**Table 6**     Performance of SAA and Algorithm 3 in W-System under the Sampled Distribution

| Sample Size | n | SAA | | Algorithm 3 | | Step | Relative Difference in Capacity |
|---|---|---|---|---|---|---|---|
| | | Capacity | Time (Seconds) | Capacity | Time (Seconds) | | (Capacity by Algorithm 3-Capacity by SAA)/Capacity by SAA |
| 50 | 5 | 137 | 0.46 | 135 | 78.28 | 1 | -1.46% |
| | 10 | 248 | 0.30 | 248 | 411.39 | 5 | 0.00% |
| | 15 | 377 | 0.37 | 375 | 1557.96 | 6 | -0.53% |
| | 20 | 490 | 0.36 | 486 | 934.92 | 4 | -0.82% |
| | 25 | 624 | 1.24 | 618 | 1678.33 | 7 | -0.96% |
| | 30 | 740 | 0.66 | 737 | 1629.37 | 7 | -0.41% |
| | 35 | 869 | 1581.35 | 861 | 1580.79 | 7 | -0.92% |
| | 40 | 977 | 1.05 | 977 | 2701.97 | 12 | 0.00% |
| | 45 | 1127 | 0.75 | 1130 | 4764.06 | 22 | 0.27% |
| | 50 | 1229 | 5.66 | 1228 | 3141.24 | 15 | -0.08% |
| 100 | 5 | 133 | 0.31 | 134 | 701.59 | 5 | 0.75% |
| | 10 | 245 | 1.02 | 248 | 1743.75 | 7 | 1.22% |
| | 15 | 382 | 2.76 | 386 | 2635.15 | 12 | 1.05% |
| | 20 | 493 | 59.63 | 502 | 3466.15 | 17 | 1.83% |
| | 25 | 615 | 3.09 | 625 | 2844.83 | 15 | 1.63% |
| | 30 | 736 | 0.91 | 743 | 2063.59 | 10 | 0.95% |
| | 35 | 859 | 79.44 | 865 | 2262.32 | 13 | 0.70% |
| | 40 | 964 | 2.15 | 980 | 3014.55 | 17 | 1.66% |
| | 45 | 1109 | 6.44 | 1129 | 4327.14 | 24 | 1.80% |
| | 50 | 1208 | 1.76 | 1233 | 3986.90 | 22 | 2.07% |
| 500 | 5 | 136 | 1.24 | 137 | 1011.98 | 4 | 0.74% |
| | 10 | 250 | 8.85 | 251 | 1170.64 | 6 | 0.40% |
| | 15 | 378 | 7.28 | 381 | 1740.91 | 9 | 0.79% |
| | 20 | 491 | 22.41 | 495 | 1937.38 | 10 | 0.81% |
| | 25 | 617 | 2478.16 | 622 | 2476.43 | 14 | 0.81% |
| | 30 | 727 | 116.04 | 736 | 3078.88 | 17 | 1.24% |
| | 35 | OOM | | 866 | 2355.64 | 19 | |
| | 40 | OOM | | 970 | 2534.21 | 22 | |
| | 45 | OOM | | 1128 | 3471.02 | 25 | |
| | 50 | OOM | | 1231 | 3561.52 | 23 | |

To ensure the capacity is also feasible for the original distribution, we implement Algorithm 3 and the SA heuristic for each problem instance using the original (normal) distribution. The results are presented in Table 8. We observe that the total capacity obtained by directly implementing the SA heuristic is always much higher than the total capacity obtained by Algorithm 3. The difference is between 29.8% and 35.7% for the W-system, and between 65% and 72.9% for the non-W-system.

We demonstrate in this subsection how the Max-Weighted-Service policy can be used in a simple local search algorithm to compute a near-optimal capacity level. We leave it for future research to develop more sophisticated and computationally efficient algorithms with better performance.

## 6.  Conclusions

In this paper, we present a general framework to study two-stage capacity allocation and demand fulfillment with individual service constraints. We propose the Max-Weighted-Service policy and prove its asymptotic optimality for a general class of problems. When the set of feasible fulfilled demand is a polymatroid and when both the allocation cost function and the service measure function are linear in fulfilled demand, a randomized index policy is asymptotically optimal. Moreover, we formulate our model as a minimax stochastic program so that the optimal capacity level can

**Table 7** Performance of SAA and Algorithm 3 in Non-W-System under the Sampled Distribution

| Sample Size | n | SAA Capacity | SAA Time (Seconds) | Algorithm 3 Capacity | Algorithm 3 Time (Seconds) | Step | Relative Difference in Capacity (Capacity by Algorithm 3-Capacity by SAA)/Capacity by SAA |
|---|---|---|---|---|---|---|---|
| | 5 | 217 | 0.87 | 214 | 2043.32 | 5 | -1.38% |
| | 10 | 401 | 1.23 | 401 | 5222.18 | 12 | 0.00% |
| | 15 | 615 | 1.45 | 629 | 9538.34 | 25 | 2.28% |
| | 20 | 810 | 3.85 | 804 | 5074.73 | 18 | -0.74% |
| | 25 | 810 | 3.85 | 804 | 5074.73 | 18 | -0.74% |
| 50 | 30 | 1241 | 16.29 | 1257 | 14537.71 | 44 | 1.29% |
| | 35 | 1461 | 30.18 | 1469 | 15358.62 | 49 | 0.55% |
| | 40 | 1641 | 23.74 | 1647 | 13937.40 | 54 | 0.37% |
| | 45 | 1906 | 72.80 | 1905 | 15189.89 | 54 | -0.05% |
| | 50 | 2076 | 65.31 | 2108 | 23989.59 | 85 | 1.54% |
| | 5 | 210 | 0.96 | 212 | 2506.92 | 7 | 0.95% |
| | 10 | 395 | 1.73 | 403 | 4504.88 | 13 | 2.03% |
| | 15 | 630 | 14.72 | 644 | 8662.06 | 28 | 2.22% |
| | 20 | 821 | 1286.46 | 834 | 9627.64 | 32 | 1.58% |
| | 25 | 1035 | 2140.52 | 1061 | 16142.44 | 47 | 2.51% |
| 100 | 30 | 1238 | 5504.34 | 1246 | 10230.32 | 33 | 0.65% |
| | 35 | 1448 | 4728.77 | 1458 | 12012.32 | 39 | 0.69% |
| | 40 | 1629 | 14281.35 | 1658 | 14278.87 | 58 | 1.78% |
| | 45 | 1892 | 17240.81 | 1914 | 17237.81 | 60 | 1.16% |
| | 50 | 2065 | 25316.04 | 2098 | 25314.95 | 73 | 1.60% |
| | 5 | 216 | 9.53 | 221 | 4366.93 | 11 | 2.31% |
| | 10 | 407 | 32.80 | 408 | 3411.44 | 12 | 0.25% |
| | 15 | 630 | 6651.68 | 634 | 6650.38 | 23 | 0.63% |
| | 20 | 819 | 6868.56 | 823 | 6866.42 | 24 | 0.49% |
| | 25 | 1039 | 11074.69 | 1047 | 11071.73 | 35 | 0.77% |
| 500 | 30 | 1241 | 11162.33 | 1251 | 11157.59 | 32 | 0.81% |
| | 35 | OOM | | 1463 | 12377.60 | 43 | |
| | 40 | OOM | | 1648 | 15255.82 | 53 | |
| | 45 | OOM | | 1932 | 17654.68 | 75 | |
| | 50 | OOM | | 2118 | 21209.06 | 88 | |

**Table 8** Capacity Solved by Algorithm 3 and SA Heuristic under the Original Distribution

| n | W-System Algorithm 3 | W-System SA Heuristic | Non-W-System Algorithm 3 | Non-W-System SA Heuristic |
|---|---|---|---|---|
| 5 | 137 | 186 | 220 | 363 |
| 10 | 251 | 337 | 411 | 711 |
| 15 | 380 | 501 | 636 | 1075 |
| 20 | 494 | 654 | 836 | 1427 |
| 25 | 627 | 817 | 1075 | 1794 |
| 30 | 741 | 969 | 1262 | 2140 |
| 35 | 868 | 1133 | 1489 | 2506 |
| 40 | 970 | 1279 | 1655 | 2842 |
| 45 | 1124 | 1459 | 1955 | 3239 |
| 50 | 1229 | 1603 | 2129 | 3571 |

be computed or approximated by applying existing first-order optimization algorithms, such as the mirror descent SA algorithm.

Our results are derived by a strong duality result formulated in Theorem 2, for any fixed capacity level. We demonstrate the potential use of strong duality to analyze optimal capacity level for problems studied in the literature, i.e., inventory pooling and assemble to order with two products and with i.i.d. uniform distribution.

Theorem 2 can also be used to obtain further analytical results. For example, consider inventory pooling with i.i.d. common strictly increasing continuous demand distribution function, equal target service level $\beta$, and Type I service constraints. When the number of customers goes to infinity, we are able to derive a closed-form expression for the asymptotically optimal inventory level per customer. We call a per customer capacity level $c^*$ asymptotically optimal if $c^*$ is the smallest value of $c$ such that the capacity level $n \cdot c$ is feasible as $n \to \infty$. Then we have

$$c^* = \max_{\xi}\{\xi - \frac{1}{\beta}\,\mathrm{E}_{\tilde{D}}[(\xi - \tilde{D})^+]\}.$$

When the common distribution has a finite mean $\mu$ and a finite standard deviation $\sigma$, we are able to obtain the following bounds on $c^*$

$$\beta\mu - \sqrt{\beta(1-\beta)}\sigma \leq c^* \leq \beta\mu.$$

The formal proof is relegated to Appendix I.

In our formulation (10), there is exactly one service constraint for each customer. However, our approach continues to apply even when there are multiple service constraints per customer. For example, a customer can impose both Type I and Type II service constraints. Assume that service constraints are given by

$$\mathrm{E}_{\tilde{\boldsymbol{\phi}},\tilde{\mathbf{D}}}[R_j^k(s_j(\tilde{\boldsymbol{\phi}},\mathbf{c},\tilde{\mathbf{D}}),\tilde{D}_j)] \geq \beta_j^k, \quad \forall j \in \mathcal{N}, k = 1, ..., K,$$

where $R_j^k$ is the $k$th service measure function of customer $j$. Then we can define

$$\hat{g}(\mathbf{w},\mathbf{c},\mathbf{D}) = \max \sum_{j \in \mathcal{N}} \sum_{k=1}^{K} w_j^k R_j^k(s_j, D_j)$$
$$\text{s.t. } \mathbf{s} \in P(\mathbf{c},\mathbf{D})$$

Then we obtain as generalization of Corollary 1. More specifically, $\mathbf{c}$ is feasible if and only if

$$\mathrm{E}_{\mathbf{D}}[\hat{g}(\mathbf{w},\mathbf{c},\mathbf{D})] \geq \sum_{j \in \mathcal{N}} \sum_{k=1}^{K} w_j^k \beta_j^k \quad \forall \mathbf{w} \geq 0.$$

We can also show that a generalization of the Max-Weighted-Service policy is optimal.

Instead of minimizing the total capacity level subject to service level constraints, we can maximize a (concave) function of achieved service levels for any given fixed capacity level. The new formulation can always be formulated as a concave-convex stochastic saddle point problem, and thus can be solved by existing stochastic approximation algorithms.

We have discussed that our capacity rationing policy is applicable to a periodic-review infinite time horizon model where we assume the capacity is perishable and unmet demand is lost. It is

more challenging to analyze the problem when unmet demand is backlogged; see Shi et al. (2019). This is a topic of our ongoing research.

## Acknowledgments

## References

Agrawal, Narendra, Morris A Cohen. 2001. Optimal material control in an assembly system with component commonality. *Naval Research Logistics (NRL)* **48**(5) 409–429.

Alptekinoğlu, Aydın, Arunava Banerjee, Anand Paul, Nikhil Jain. 2013. Inventory pooling to deliver differentiated service. *Manufacturing & Service Operations Management* **15**(1) 33–44.

Anupindi, Ravi, Yehuda Bassok, Eitan Zemel. 2001. A general framework for the study of decentralized distribution systems. *Manufacturing & Service Operations Management* **3**(4) 349–368.

Asadpour, Arash, Xuan Wang, Jiawei Zhang. 2020. Online resource allocation with limited flexibility. *Management Science* **66**(2) 642–666.

Azuma, Kazuoki. 1967. Weighted sums of certain dependent random variables. *Tohoku Mathematical Journal, Second Series* **19**(3) 357–367.

Baker, Kenneth R. 1985. Safety stocks and component commonality. *Journal of Operations Management* **6**(1) 13–22.

Baker, Kenneth R, Michael J Magazine, Henry LW Nuttle. 1986. The effect of commonality on safety stock in a simple inventory model. *Management Science* **32**(8) 982–988.

Bassamboo, Achal, Ramandeep S Randhawa, Jan A Van Mieghem. 2010. Optimal flexibility configurations in newsvendor networks: Going beyond chaining and pairing. *Management Science* **56**(8) 1285–1303.

Bassok, Yehuda, Ravi Anupindi, Ram Akella. 1999. Single-period multiproduct inventory models with substitution. *Operations Research* **47**(4) 632–642.

Blackwell, David. 1956. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics* **6**(1) 1–8.

Borel, M Émile. 1909. Les probabilités dénombrables et leurs applications arithmétiques. *Rendiconti del Circolo Matematico di Palermo (1884-1940)* **27**(1) 247–271.

Bruss, F Thomas, James B Robertson. 1991. 'wald's lemma'for sums of order statistics of iid random variables. *Advances in applied probability* **23**(3) 612–623.

Cachon, Gérard, Christian Terwiesch. 2008. *Matching supply with demand*. McGraw-Hill Publishing.

Chen, Xi, Jiawei Zhang, Yuan Zhou. 2015. Optimal sparse designs for process flexibility via probabilistic expanders. *Operations Research* **63**(5) 1159–1176.

Chou, Mabel C, Geoffrey A Chua, Chung-Piaw Teo, Huan Zheng. 2010. Design for process flexibility: Efficiency of the long chain and sparse structure. *Operations research* **58**(1) 43–58.

Chou, Mabel C, Geoffrey A Chua, Chung-Piaw Teo, Huan Zheng. 2011. Process flexibility revisited: The graph expander and its applications. *Operations Research* **59**(5) 1090–1105.

Désir, Antoine, Vineet Goyal, Yehua Wei, Jiawei Zhang. 2016. Sparse process flexibility designs: is the long chain really optimal? *Operations Research* **64**(2) 416–431.

Eppen, Gary D. 1979. Note—effects of centralization on expected costs in a multi-location newsboy problem. *Management science* **25**(5) 498–501.

Gerchak, Yigal, Mordechai Henig. 1989. Component commonality in assemble-to-order systems: Models and properties. *Naval Research Logistics (NRL)* **36**(1) 61–68.

Gerchak, Yigal, Michael J Magazine, A Bruce Gamble. 1988. Component commonality with service level requirements. *Management science* **34**(6) 753–760.

Gurvich, Itai, James Luedtke, Tolga Tezcan. 2010. Staffing call centers with uncertain demand forecasts: A chance-constrained optimization approach. *Management Science* **56**(7) 1093–1115.

Hazan, Elad. 2019. Introduction to online convex optimization. *arXiv preprint arXiv:1909.05207* .

He, Simai, Jiawei Zhang, Shuzhong Zhang. 2012. Polymatroid optimization, submodularity, and joint replenishment games. *Operations research* **60**(1) 128–137.

Hou, I-H, Vivek Borkar, PR Kumar. 2009. A theory of qos for wireless. *Proc. IEEE Conf. Comput. Commun*. IEEE, 486–494.

Jordan, William C, Stephen C Graves. 1995. Principles on the benefits of manufacturing process flexibility. *Management Science* **41**(4) 577–594.

Juditsky, Anatoli, Arkadi Nemirovski. 2011. First order methods for nonsmooth convex large-scale optimization, i: general purpose methods. *Optimization for Machine Learning* 121–148.

Lee, Hau L. 1996. Effective inventory and service management through product and process redesign. *Operations Research* **44**(1) 151–159.

Liu, Xiao, Simge Küçükyavuz, James Luedtke. 2016. Decomposition algorithms for two-stage chance-constrained programs. *Mathematical Programming* **157**(1) 219–243.

Lyu, Guodong, Wang-Chi Cheung, Mabel C Chou, Chung-Piaw Teo, Zhichao Zheng, Yuanguang Zhong. 2019. Capacity allocation in flexible production networks: Theory and applications. *Management Science* **65**(11) 5091–5109.

Lyu, Guodong, Mabel Chou, Chung-Piaw Teo, Zhichao Zheng, Yuanguang Zhong. 2017. Capacity allocation with stock-out probability targets: Theory and applications .

Martin, Kipp, Christopher Thomas Ryan, Matt Stern. 2016. The slater conundrum: duality and pricing in infinite-dimensional optimization. *SIAM Journal on Optimization* **26**(1) 111–138.

Megiddo, Nimrod. 1974. Optimal flows in networks with multiple sources and sinks. *Mathematical Programming* **7**(1) 97–107.

Mieghem, Jan A Van, Nils Rudi. 2002. Newsvendor networks: Inventory management and capacity investment with discretionary activities. *Manufacturing & Service Operations Management* **4**(4) 313–335.

Mirchandani, Prakash, Ajay K Mishra. 2002. Component commonality: Models with product-specific service constraints. *Production and Operations Management* **11**(2) 199–215.

Nemirovski, Arkadi, Anatoli Juditsky, Guanghui Lan, Alexander Shapiro. 2009. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization* **19**(4) 1574–1609.

Shapiro, Alexander. 2001. On duality theory of conic linear problems. *Semi-infinite programming*. Springer, 135–165.

Shi, Cong, Yehua Wei, Yuan Zhong. 2019. Process flexibility for multiperiod production systems. *Operations Research* **67**(5) 1300–1320.

Simchi-Levi, David, Yehua Wei. 2012. Understanding the performance of the long chain and sparse designs in process flexibility. *Operations research* **60**(5) 1125–1141.

Simchi-Levi, David, Yehua Wei. 2015. Worst-case analysis of process flexibility designs. *Operations Research* **63**(1) 166–185.

Sion, Maurice. 1958. On general minimax theorems. *Pacific Journal of mathematics* **8**(1) 171–176.

Swaminathan, Jayashankar M, Ramesh Srinivasan. 1999. Managing individual customer service constraints under stochastic demand. *Operations Research Letters* **24**(3) 115–125.

Van Mieghem, Jan A. 1998. Investment strategies for flexible resources. *Management Science* **44**(8) 1071–1078.

Wang, Xuan, Jiawei Zhang. 2015. Process flexibility: A distribution-free bound on the performance of k-chain. *Operations Research* **63**(3) 555–571.

Welsh, Dominic JA. 2010. *Matroid theory*. Courier Corporation.

Zhang, Alex X. 1997. Demand fulfillment rates in an assembleto-order system with multiple products and dependent demands. *Production and Operations Management* **6**(3) 309–324.

Zhong, Yuanguang, Zhichao Zheng, Mabel C Chou, Chung-Piaw Teo. 2017. Resource pooling and allocation policies to deliver differentiated service. *Management Science* **64**(4) 1555–1573.

## Appendix A: Proof of Lemma 1

*Proof:* The equations hold due to our definition of deterministic policies. To see this, consider any given deterministic policy $\phi \in \Phi$. It determines a unique allocation $(\mathbf{y}(\phi, \mathbf{c}, \mathbf{D}), \mathbf{s}(\phi, \mathbf{c}, \mathbf{D})) \in$

$P(\mathbf{c}, \mathbf{D})$ for every demand realization $\mathbf{D}$. Thus $(\mathbf{y}(\phi, \mathbf{c}, \mathbf{D}), \mathbf{s}(\phi, \mathbf{c}, \mathbf{D}))$ is a feasible solution to (16). This implies that

$$f(\mathbf{y}(\phi, \mathbf{c}, \mathbf{D})) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(s_j(\phi, \mathbf{c}, \mathbf{D}), D_j) \geq g(\mathbf{w}, \mathbf{c}; \mathbf{D})$$

and thus

$$F(\mathbf{w}, \phi) = E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] - \sum_{j \in \mathcal{N}} w_j \cdot E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \geq E_{\tilde{\mathbf{D}}}[g(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})]$$

holds for any $\phi \in \Phi$. Thus,

$$\inf_{\phi \in \Phi} F(\mathbf{w}, \phi) \geq E_{\tilde{\mathbf{D}}}[g(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})]$$

On the other hand, by the definition of $\phi_{\mathbf{w}}$, we have that

$$g(\mathbf{w}, \mathbf{c}; \mathbf{D}) = f(\mathbf{y}(\phi_{\mathbf{w}}, \mathbf{c}, \mathbf{D})) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(s_j(\phi_{\mathbf{w}}, \mathbf{c}, \mathbf{D}), D_j)$$

and thus

$$E_{\tilde{\mathbf{D}}}[g(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})] = F(\mathbf{w}, \phi_{\mathbf{w}}) \geq \inf_{\phi \in \Phi} F(\mathbf{w}, \phi)$$

Therefore, equality (17) holds.    $\square$

## Appendix B:  Proof of Theorem 2

We first prove the following lemma, which will lead to our final result. Its proof is relegated to the Appendix C.

LEMMA 3. *Under Assumption 2, there exists a subset of deterministic policies $\Phi_W$ such that for any $\mathbf{w} \geq 0$, we have*

$$\inf_{\phi \in \Phi} F(\mathbf{w}, \phi) = F(\mathbf{w}, \phi_{\mathbf{w}}) = \min_{\phi \in \Phi_W} F(\mathbf{w}, \phi). \tag{36}$$

*Moreover, for any sequence of randomized policies $\{\lambda_k\}_{k \geq 1}$ such that $\lambda_k \in \chi_W := \{\lambda \geq 0 : \int_{\phi \in \Phi_W} d\lambda(\phi) = 1\}$ for each integer $k$, there exists a policy $\hat{\lambda} \in \chi$ such that*

$$\int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] d\hat{\lambda}(\phi) \leq \limsup_{k \to \infty} \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\phi) \tag{37}$$

*and*

$$\int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\hat{\lambda}(\phi) \geq \liminf_{k \to \infty} \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_k(\phi), \quad \forall j \in \mathcal{N}. \tag{38}$$

We are now ready to prove the strong duality between (11) and (12). We first use weak duality to show that Obj (11) $\geq$ Obj (12). In order to prove the other direction, we construct a sequence of randomized policy $\lambda_k$ such that the achieved service level is at least $\beta_j - \frac{\hat{C}}{k}$ for each $j \in \mathcal{N}$, for some constant $\hat{C} > 0$. Then, by letting $k \to \infty$, we use Lemma 3 to show that there exists a randomized policy such that the target service level $\beta_j$ is achieved for each $j \in \mathcal{N}$ and the expected allocation cost is upper bounded by Obj (12), which completes our proof.

*Proof of Theorem 2:*    We denote Obj (11) as the objective value of (11) and denote Obj (12) as the objective value of (12). When (11) is feasible, we get

$$\text{Obj (11)} \geq \text{Obj (12)} \tag{39}$$

by weak duality (Shapiro, 2001). In the remaining part of the proof, we assume that Obj (12) is finite, and we prove that (11) is feasible and Obj (11) = Obj (12).

Note that

$$\inf_{\lambda \in \chi} L(\mathbf{w}, \lambda) = \sum_{j \in \mathcal{N}} w_j \cdot \beta_j + \inf_{\phi \in \Phi} F(\mathbf{w}, \phi) = \sum_{j \in \mathcal{N}} w_j \cdot \beta_j + \inf_{\phi \in \Phi_W} F(\mathbf{w}, \phi) = \inf_{\lambda \in \chi_W} L(\mathbf{w}, \lambda) \tag{40}$$

where the second equality holds due to Lemma 3. Then we have

$$\sup_{\mathbf{w} \geq 0} \inf_{\lambda \in \chi_W} L(\mathbf{w}, \lambda) = \text{Obj (12)}$$

For each integer $k > 1$, we define the set $W_k = \{\mathbf{w} \geq 0 : \sum_{j \in \mathcal{N}} w_j \leq k\}$. Obviously, it holds that

$$\sup_{\mathbf{w} \in W_k} \inf_{\lambda \in \chi_W} L(\mathbf{w}, \lambda) \leq \sup_{\mathbf{w} \geq 0} \inf_{\lambda \in \chi_W} L(\mathbf{w}, \lambda) = \text{Obj (12)}$$

By definition, $\chi_W$ is a convex set. Moreover, note that $W_k$ is a convex compact set and $L(\mathbf{w}, \lambda)$ is linear in $\mathbf{w}$ (resp. $\lambda$) when $\lambda$ (resp. $\mathbf{w}$) is fixed. Then by Sion's minimax theorem (Sion, 1958), we must have

$$\inf_{\lambda \in \chi_W} \sup_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda) = \inf_{\lambda \in \chi_W} \max_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda) = \max_{\mathbf{w} \in W_k} \inf_{\lambda \in \chi_W} L(\mathbf{w}, \lambda) = \sup_{\mathbf{w} \in W_k} \inf_{\lambda \in \chi_W} L(\mathbf{w}, \lambda) \leq \text{Obj (12)}$$

Thus, there exists a randomized policy $\lambda_k \in \chi_W$ such that

$$\sup_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda_k) \leq \inf_{\lambda \in \chi_W} \sup_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda) + \frac{1}{k} \leq \text{Obj (12)} + \frac{1}{k} \tag{41}$$

Denote $\hat{C} = \text{Obj (12)} + 1$. We now claim that $\lambda_k$ achieves a service level at least $\beta_j - \frac{\hat{C}}{k}$ for each $j \in \mathcal{N}$, i.e.

$$\int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_l(\phi) \geq \beta_j - \frac{\hat{C}}{k}. \tag{42}$$

Otherwise, suppose there exists a $j_0 \in \mathcal{N}$ such that

$$\hat{C} < k \cdot \left( \beta_{j_0} - \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_{j_0}(s_{j_0}(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_{j_0})] d\lambda_k(\phi) \right).$$

Then we define $\hat{\mathbf{w}} \in W_k$ such that $\hat{w}_{j_0} = k$ and $\hat{w}_j = 0$ for all $j \neq j_0$. By construction, we have

$$\begin{aligned}
\hat{C} &< \hat{w}_{j_0} \cdot \left( \beta_{j_0} - \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_{j_0}(s_{j_0}(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_{j_0})] d\lambda_k(\phi) \right) \\
&= \sum_{j \in N} \hat{w}_j \cdot \left( \beta_j - \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_k(\phi) \right) \\
&\leq \sum_{j \in N} \hat{w}_j \cdot \left( \beta_j - \int_{\phi \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_k(\phi) \right) + E_{\tilde{\mathbf{D}}}[f(y(\phi, \mathbf{c}, \tilde{\mathbf{D}}))] \\
&= L(\hat{\mathbf{w}}, \lambda_k) \leq \sup_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda_k) \leq \hat{C}
\end{aligned}$$

where the second inequality holds since the allocation cost function $f$ is non-negative and the last inequality follows from (41) since $1 \geq \frac{1}{k}$. This is a contradiction.

From Lemma 3, for the sequence $\{\lambda_k\}_{k \geq 1}$, there exists a randomized policy $\hat{\lambda} \in \chi$ such that (37) and (38) hold. Specifically, for each $j \in \mathcal{N}$, we have

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\hat{\lambda}(\boldsymbol{\phi}) \geq \liminf_{k \to \infty} \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_k(\boldsymbol{\phi})$$

$$\geq \liminf_{k \to \infty} \beta_j - \frac{\hat{C}}{k} = \beta_j$$

Thus, (11) is feasible and $\hat{\lambda}$ is a feasible solution to (11). Moreover, from (37), we have

$$\text{Obj (11)} \leq \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\hat{\lambda}(\boldsymbol{\phi}) \leq \limsup_{k \to \infty} \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\boldsymbol{\phi})$$

$$\leq \limsup_{k \to \infty} \sup_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda_k) \leq \limsup_{k \to \infty} \text{Obj (12)} + \frac{1}{k} = \text{Obj (12)}$$

where the third inequality follows from the fact that

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\boldsymbol{\phi}) = L(\mathbf{0}, \lambda_k) \leq \sup_{\mathbf{w} \in W_k} L(\mathbf{w}, \lambda_k), \quad \forall k > 1$$

and the last inequality follows from (41). Together with the weak duality established in (39), we have Obj (11) = Obj (12). $\square$

### Appendix C: Proof of Lemma 3

We first present the following well-known result, which will be useful for our proof of Lemma 3.

THEOREM 6. *(Banach-Saks theorem) Let $\{\mathbf{x}^k\}_{k=1}^{\infty}$ be a bounded sequence in the Hilbert space $\mathcal{H}$, then there exists a subsequence $\{n_k\}_{k=1}^{\infty}$ of $\{1, 2, \ldots, k, \ldots\}$ and a point $\mathbf{x} \in \mathcal{H}$ such that*

$$\frac{1}{k} \cdot \sum_{l=1}^{k} \mathbf{x}^{n_l}$$

*converges strongly to $\mathbf{x}$ as $k \to \infty$.*

Now we are ready to prove Lemma 3. The following is an outline of the proof of Lemma 3. We first represent each policy as an element in a Hilbert space and interpret each term in (37) and (38) as an inner product that defines the metric of the Hilbert space. We then apply the Banach-Saks theorem to prove weak convergence of the running average of a subsequence to a limiting element in the Hilbert space, which directly implies that the cost and the service level of the running average converge to those of the limiting element. Finally, we show that the limiting element can be interpreted back as a policy, which establishes the existence of $\hat{\lambda}$ in Lemma 3.

*Proof of Lemma 3:* Note that under Assumption 2, the service measure function for each $j \in \mathcal{N}$ is piece-wise linear in $s_j$. Thus, the objective function of (16) is a piece-wise linear function. Then, problem (16) can be solved by first enumerating all possible vectors $\mathbf{k} \in K := \{(k_1, k_2, \cdots, k_n) : k_j = 1, 2, \ldots, K_j\}$, such that

$$a_{j,k_j} \cdot D_j \le s_j < a_{j,k_j+1} \cdot D_j,$$

and then for each set of values $\mathbf{k}$ solving

$$
\begin{aligned}
g(\mathbf{w}, \mathbf{c}; \mathbf{D}; \mathbf{k}) := \min_{\mathbf{s}, \mathbf{y}} \ & f(y) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(s_j, D_j) \\
\text{s.t.} \quad & (\mathbf{s}, \mathbf{y}) \in P(\mathbf{c}, \mathbf{D}) \\
& a_{j,k_j} \cdot D_j \le s_j \le a_{j,k_j+1} \cdot D_j \ \ \forall j \in \mathcal{N}.
\end{aligned}
\tag{43}
$$

Note that in the problem (43), we replace the constraint $s_j < a_{j,k_j+1} \cdot D_j$ with $s_j \le a_{j,k_j+1} \cdot D_j$. This is because for fixed $k_j$, we can simply define $R_j(a_{j,k_j+1} \cdot D_j, D_j) = \lim_{s_j \to (a_{j,k_j+1} \cdot D_j)-} R_j(s_j, D_j)$ denoting the left limit when solving (43). Then the minimum value of $g(\mathbf{w}, \mathbf{c}; \mathbf{D}; \mathbf{k})$ over all possible $\mathbf{k} \in K$ is still the same as the optimal value of the original problem (16).

By definition of $P(\mathbf{c}, \mathbf{D})$, the linear program (43) can be reformulated as:

$$
\begin{aligned}
g(\mathbf{w}, \mathbf{c}, \mathbf{D}; \mathbf{k}) = \min \ & \mathbf{r}_{\mathbf{w}, \mathbf{D}}^T \hat{\mathbf{y}} \\
\text{s.t.} \quad & \mathbf{A}\hat{\mathbf{y}} = \mathbf{v_D} \\
& \hat{\mathbf{y}} \ge 0
\end{aligned}
\tag{44}
$$

by choosing the appropriate $\mathbf{r}_{\mathbf{w},\mathbf{D}}, \mathbf{A}$ and $\mathbf{v_D}$, where $\mathbf{v_D}$ is independent of $\mathbf{w}$ and $\mathbf{A}$ is independent of $\mathbf{w}$ and $\mathbf{D}$. Although all the coefficients $\mathbf{r}_{\mathbf{w},\mathbf{D}}, \mathbf{v_D}$ and $\mathbf{A}$ should also be dependent on $\mathbf{k}$, we drop the dependency for simplicity of notation. The dependence of $\mathbf{v_D}$ on $\mathbf{c}$ is also dropped since $\mathbf{c}$ is fixed. We now focus on LP (44).

Denote $\mathcal{D}$ as the support of the demand distribution. Given $\mathbf{k}$, for each $\mathbf{w} \ge 0$ and $\mathbf{D} \in \mathcal{D}$, there could be multiple optimal solutions. However, it is enough for us to only consider one optimal *basic solution* that is determined by a basis $\mathbf{b} \in \mathcal{B}$, where $\mathcal{B}$ is the set of all bases of $\mathbf{A}$. Since $\mathbf{A}$ has a finite size, the total number of all bases, $|\mathcal{B}|$ should be finite. Then for any $\mathbf{w} \ge 0$ and $\mathbf{D} \in \mathcal{D}$, an optimal solution of (16) is uniquely determined by an element in the finite set $\mathcal{V} = \mathcal{K} \times \mathcal{B}$. For the rest of the proof, we only consider such optimal solutions. Without loss of generality, we sort the elements in the set $\mathcal{V}$ in a fixed sequence and we will use the order of an element in this sequence to denote this element by abuse of notation.

For each fixed $\mathbf{w} \in W$, we define a deterministic policy $\hat{\phi}_{\mathbf{w}}$ as follows. For any $\mathbf{D} \in \mathcal{D}$, let $\hat{\phi}_{\mathbf{w}}(\mathbf{c}, \mathbf{D})$ be the specified optimal solution of (16) determined by one element from $\mathcal{V}$. As a direct consequence of this definition, for each $\mathbf{w} \in W$, we have that

$$\hat{\phi}_{\mathbf{w}} \in \operatorname{argmin}_{\phi \in \Phi} F(\mathbf{w}, \phi).$$

We define $\Phi_W = \{\hat{\boldsymbol{\phi}}_{\mathbf{w}} : \forall \mathbf{w} \in W\}$. Then, our proof of (36) is finished. In the remaining part of the proof, we prove (37) and (38).

For each $\mathbf{D} \in \mathcal{D}$ and each element $v \in \mathcal{V}$, we further denote $a(v, \mathbf{D})$ as the basic solution of (44) determined by the element $v$ if feasible, i.e., $a(v, \mathbf{D}) \in P(\mathbf{c}, \mathbf{D})$. If infeasible, we simply denote $a(v, \mathbf{D}) = \mathbf{0}$. Then, from definition, for each $\mathbf{w} \in W$ and each $\mathbf{D} \in \mathcal{D}$, the allocation $\hat{\boldsymbol{\phi}}_{\mathbf{w}}(\mathbf{c}, \mathbf{D}) \in \{a(v, \mathbf{D})\}_{\forall v \in \mathcal{V}}$.

We now focus on the sequence of randomized policies $\{\lambda_k\}$ such that $\lambda_k \in \chi_W := \{\lambda \geq 0 : \int_{\boldsymbol{\phi} \in \Phi_W} d\lambda(\boldsymbol{\phi}) = 1\}$. From the above argument, we know that for any $k$ and any $\mathbf{D} \in \mathcal{D}$, the allocation of $\lambda_k$ is simply a randomization over $\{a(v, \mathbf{D})\}_{\forall v \in \mathcal{V}}$. Then, we define a vector $\boldsymbol{\psi}^k(\mathbf{D}) \in \mathbb{R}^{|\mathcal{V}|}$, where $|\mathcal{V}|$ denotes the cardinality of the finite set $\mathcal{V}$, such that the $v$-th component of $\boldsymbol{\psi}^k(\mathbf{D})$, denoted as $\psi_v^k(\mathbf{D})$, denotes the probability that the allocation of $\lambda_k$ equals $a(v, \mathbf{D})$, given demand realization $\mathbf{D}$. Obviously, we have that

$$\boldsymbol{\psi}^k(\mathbf{D}) \in \mathcal{L} := \{\mathbf{x} \in \mathbb{R}^{|\mathcal{V}|} : \mathbf{x} \geq 0, \ \sum_{v=1}^{|\mathcal{V}|} x_v = 1\}, \quad \forall \mathbf{D} \in \mathcal{D} \tag{45}$$

Following this definition, each randomized policy $\lambda_k$ is equivalently represented by $\boldsymbol{\psi}^k = (\boldsymbol{\psi}(\mathbf{D}), \forall \mathbf{D} \in \mathcal{D})$, which is a measurable function mapping the set $\mathcal{D} \subset \mathbb{R}^n$ to the set $\mathcal{L} \subset \mathbb{R}^{|\mathcal{V}|}$. From (45), it is easy to see that

$$\|\psi_v^k\|_{L^2}^2 = \int_{\mathbf{D} \in \mathcal{D}} |\psi_v^k(\mathbf{D})|^2 d\mu(\mathbf{D}) \leq 1, \ \ \forall v = 1, \ldots, |\mathcal{V}| \tag{46}$$

Here, $\mu$ denotes the measure over $\mathcal{D}$ specified by the demand distribution and $\|\cdot\|_{L^2}$ denotes the $L^2$-norm. Then, for each $k$ and each $v = 1, \ldots, |\mathcal{V}|$, we conclude that $\psi_v^k \in L^2(\mathcal{D}, \mu)$, where $L^2(\mathcal{D}, \mu)$ denotes the $L^2$ space containing all measurable functions over the set $\mathcal{D}$ with finite $L^2$-norm.

For each $v = 1, \ldots, |\mathcal{V}|$, we further denote $L_v$ as a copy of the space $L^2(\mathcal{D}, \mu)$, which is a Hilbert space. We then denote $\mathcal{H}$ as the direct sum of the spaces $\{L_v\}_{v=1}^{|\mathcal{V}|}$, i.e.,

$$\mathcal{H} = \bigoplus_{v=1}^{|\mathcal{V}|} L_v := \{\forall \boldsymbol{\varphi} = (\varphi_v, v = 1, \ldots, |\mathcal{V}|) \text{ such that } \varphi_v \in L_v \text{ for each } v\}$$

Clearly, $\mathcal{H}$ is still a Hilbert space, equipped with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ defined as:

$$\langle \boldsymbol{\varphi}^1, \boldsymbol{\varphi}^2 \rangle_{\mathcal{H}} = \sum_{v=1}^{|\mathcal{V}|} \int_{\mathbf{D} \in \mathcal{D}} \varphi_v^1(\mathbf{D}) \cdot \varphi_v^2(\mathbf{D}) d\mu(\mathbf{D}), \quad \forall \boldsymbol{\varphi}^1, \boldsymbol{\varphi}^2 \in \mathcal{H}$$

Denote $\|\cdot\|_{\mathcal{H}}$ as the norm on the Hilber space $\mathcal{H}$ induced by the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$. For each $k$, it holds directly that $\boldsymbol{\psi}^k \in \mathcal{H}$ and from (46), we have

$$\|\boldsymbol{\psi}^k\|_{\mathcal{H}}^2 = \sum_{v=1}^{|\mathcal{V}|} \int_{\mathbf{D} \in \mathcal{D}} |\psi_v^k(\mathbf{D})|^2 d\mu(\mathbf{D}) \leq |\mathcal{V}| \tag{47}$$

We then show that the expected allocation cost and service level of the policy $\lambda_k$ can be expressed as the inner product in the space $\mathcal{H}$. For each $v \in \mathcal{V}$ and each $\mathbf{D} \in \mathcal{D}$, we specify allocation $\mathbf{y}(v, \mathbf{D})$ and fulfillment $\mathbf{s}(v, \mathbf{D})$ such that $a(v, \mathbf{D}) = (\mathbf{y}(v, \mathbf{D}), \mathbf{s}(v, D))$. Then, we define

$$\boldsymbol{\eta}^{(0)} = (\boldsymbol{\eta}^{(0)}(\mathbf{D}), \forall \mathbf{D} \in \mathcal{D}) \quad \text{where} \quad \boldsymbol{\eta}^{(0)}(\mathbf{D}) = (f(\mathbf{y}(v, \mathbf{D})), \forall v \in \mathcal{V}) \in \mathbb{R}^{|\mathcal{V}|}$$

Note that

$$\|\boldsymbol{\eta}^{(0)}\|_{\mathcal{H}}^2 = \sum_{v=1}^{|\mathcal{V}|} \int_{\mathbf{D} \in \mathcal{D}} |(f(\mathbf{y}(v, \mathbf{D}))|^2 d\mu(\mathbf{D}) \leq \hat{C}_1 \cdot \int_{\mathbf{D} \in \mathcal{D}} \|\mathbf{D}\|_2^2 d\mu(\mathbf{D}) < \infty$$

for some constant $\hat{C}_1 > 0$, where the first inequality follows from Assumption 2b and the fact that $f(\cdot)$ is a linear function, and the second inequality follows from the demand distribution has a bounded second moment. We conclude that $\eta_v^{(0)} \in L^2(\mathcal{D}, \mu)$ for each $v = 1, \ldots, |\mathcal{V}|$ and thus $\boldsymbol{\eta}^{(0)} \in \mathcal{H}$. Moreover, since the function $f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \mathbf{D}))$ is integrable over $\mathcal{D} \times \Phi$ with respect to the measure $\mu \times \lambda_k$, then we have

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\boldsymbol{\phi}) = \int_{\boldsymbol{\phi} \in \Phi} \int_{\mathbf{D} \in \mathcal{D}} f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \mathbf{D})) d\mu(\mathbf{D}) d\lambda_k(\boldsymbol{\phi})$$
$$= \int_{\mathbf{D} \in \mathcal{D}} \int_{\boldsymbol{\phi} \in \Phi} f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \mathbf{D})) d\lambda_k(\boldsymbol{\phi}) d\mu(\mathbf{D})$$

where the second equality follows since Fubini's theorem implies that we can interchange the order of integral. Thus, it holds that

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\boldsymbol{\phi}) = \int_{\mathbf{D} \in \mathcal{D}} \sum_{v=1}^{|\mathcal{V}|} \psi_v^k(D) \cdot \eta_v^{(0)}(\mathbf{D}) d\mu(\mathbf{D}) = \langle \boldsymbol{\psi}^k, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}} \qquad (48)$$

Thus, for each $k$, the expected allocation cost of the randomized policy $\lambda_k$ can be expressed as the inner product of $\boldsymbol{\psi}^k$ and $\boldsymbol{\eta}^{(0)}$ in the space $\mathcal{H}$. Similarly, for each $j \in \mathcal{N}$, we define

$$\boldsymbol{\eta}^{(j)} = (\boldsymbol{\eta}^{(j)}(\mathbf{D}), \forall \mathbf{D} \in \mathcal{D}) \quad \text{where} \quad \boldsymbol{\eta}^{(j)}(\mathbf{D}) = (R_j(s_j(v, \mathbf{D}), D_j), \forall v \in \mathcal{V}) \in \mathbb{R}^{|\mathcal{V}|}$$

Clearly, we have that

$$\|\boldsymbol{\eta}^{(j)}\|_{\mathcal{H}}^2 = \int_{\mathbf{D} \in \mathcal{D}} \|(R_j(s_j(v, \mathbf{D}), D_j), \forall v \in \mathcal{V})\|_2^2 d\mu(\mathbf{D}) \leq C_1 \cdot \int_{\mathbf{D} \in \mathcal{D}} \|(\max\{1, D_j\}, \forall j \in \mathcal{N})\|_2^2 d\mu(\mathbf{D}) < \infty$$

Then, we conclude that for each $j \in \mathcal{N}$, $\boldsymbol{\eta}^{(j)} \in \mathcal{H}$ and we have

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\lambda_k(\boldsymbol{\phi}) = \int_{\mathbf{D} \in \mathcal{D}} \sum_{v=1}^{|\mathcal{V}|} \psi_v^k(D) \cdot \eta_v^{(j)}(\mathbf{D}) d\mu(\mathbf{D}) = \langle \boldsymbol{\psi}^k, \boldsymbol{\eta}^{(j)} \rangle_{\mathcal{H}} \qquad (49)$$

Thus, for each $k$, the service level for customer $j$ obtained by the randomized policy $\lambda_k$ can be expressed as the inner product of $\boldsymbol{\psi}^k$ and $\boldsymbol{\eta}^{(j)}$ in the space $\mathcal{H}$, for each $j \in \mathcal{N}$.

From (47), we know that the sequence $\{\boldsymbol{\psi}^k\}_{\forall k}$ is a bounded sequence in the Hilbert space $\mathcal{H}$. Then, from Banach-Saks theorem, there exists a subsequence $\{n_k\}_{\forall k}$ of $\{1, 2, \dots\}$ and $\hat{\boldsymbol{\psi}} \in \mathcal{H}$, such that the sequence $\{\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}\}_{\forall k}$ converges strongly to $\hat{\boldsymbol{\psi}}$ in the space $\mathcal{H}$. This implies that the sequence $\{\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}\}_{\forall k}$ converges weakly to $\hat{\boldsymbol{\psi}}$, then from (48), we have that

$$
\begin{aligned}
\langle \hat{\boldsymbol{\psi}}, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}} &= \lim_{k \to \infty} \langle \frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}} = \lim_{k \to \infty} \frac{1}{k} \cdot \sum_{l=1}^{k} \langle \boldsymbol{\psi}^{n_l}, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}} \leq \limsup_{k \to \infty} \langle \boldsymbol{\psi}^k, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}} \\
&= \limsup_{k \to \infty} \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\boldsymbol{\phi})
\end{aligned}
\tag{50}
$$

For each $j \in \mathcal{N}$, from (49), we have that

$$
\langle \hat{\boldsymbol{\psi}}, \boldsymbol{\eta}^{(j)} \rangle_{\mathcal{H}} = \lim_{k \to \infty} \langle \frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}, \boldsymbol{\eta}^{(j)} \rangle_{\mathcal{H}} \geq \liminf_{k \to \infty} \langle \boldsymbol{\psi}^k, \boldsymbol{\eta}^{(j)} \rangle_{\mathcal{H}} = \liminf_{k \to \infty} \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), D_j)] d\lambda_k(\boldsymbol{\phi})
\tag{51}
$$

It only remains to show that $\hat{\boldsymbol{\psi}}$ can be expressed as a randomized policy $\hat{\lambda}$. To that end, we first define a set containing all $\mathbf{D} \in \mathcal{D}$ such that $\hat{\boldsymbol{\psi}}(\mathbf{D})$ does not characterize a distribution over $\{a(v, \mathbf{D})\}_{\forall v \in \mathcal{V}}$ :

$$
\hat{\mathcal{D}} := \{\mathbf{D} \in \mathcal{D} : \hat{\boldsymbol{\psi}}(\mathbf{D}) \notin \mathcal{L}\}
$$

where the set $\mathcal{L}$ is defined in (45). We have the following result.

CLAIM 1. *It holds that $\mu(\hat{\mathcal{D}}) = 0$.*

For those $\mathbf{D} \in \hat{\mathcal{D}}$, we can simply change $\hat{\boldsymbol{\psi}}(\mathbf{D})$ into a point in the set $\mathcal{L}$. In this way, we construct a $\hat{\boldsymbol{\psi}}$ such that $\hat{\boldsymbol{\psi}}(\mathbf{D}) \in \mathcal{L}$ for any $\mathbf{D} \in \mathcal{D}$. Since $\mu(\hat{\mathcal{D}}) = 0$, we conclude that (50) and (51) still hold.

Now we show that $\hat{\boldsymbol{\psi}}$ can be characterized as a randomized policy $\hat{\lambda}$. For each $\mathbf{D} \in \mathcal{D}$, we divide the interval $[0, 1]$ into a set of sub-intervals $\{I_v(\mathbf{D})\}_{v \in \mathcal{V}}$, such that $\hat{\psi}_v(\mathbf{D}) = |I_v(\mathbf{D})|$ for each $v \in \mathcal{V}$, where $|\cdot|$ denotes the length (Lebesgue measure) of the sub-interval. Note that for each $\mathbf{D} \in \mathcal{D}$, the randomized allocation specified by $\hat{\boldsymbol{\psi}}(\mathbf{D})$ can be interpreted as picking up a point $x$ uniformly from the interval $[0, 1]$, and implementing the allocation $a(v, \mathbf{D})$ if and only if $x \in I_v(\mathbf{D})$. Thus, for each $x \in [0, 1]$, we can specify a deterministic policy

$$
\boldsymbol{\phi}_x = (\boldsymbol{\phi}_x(\mathbf{D}), \forall \mathbf{D} \in \mathcal{D}) \quad \text{where} \quad \boldsymbol{\phi}_x(\mathbf{D}) = a(v, \mathbf{D}) \text{ if and only if } x \in I_v(\mathbf{D})
$$

We define the randomized policy $\hat{\lambda}$ as the uniform distribution over the set of deterministic policies $\{\boldsymbol{\phi}_x\}_{\forall x \in [0,1]}$. We then prove (37) and (38).

For the expected allocation cost, we have that

$$
\begin{aligned}
\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\hat{\lambda}(\boldsymbol{\phi}) &= \int_{x \in [0,1]} \int_{\mathbf{D} \in \mathcal{D}} [f(\mathbf{y}(\boldsymbol{\phi}_x, \mathbf{c}, \tilde{\mathbf{D}}))] d\mu(\mathbf{D}) dx \\
&= \int_{\mathbf{D} \in \mathcal{D}} \int_{x \in [0,1]} [f(\mathbf{y}(\boldsymbol{\phi}_x, \mathbf{c}, \tilde{\mathbf{D}}))] dx \, d\mu(\mathbf{D})
\end{aligned}
$$

where the second equality follows by noting that $f(\mathbf{y}(\boldsymbol{\phi}_x, \mathbf{c}, \tilde{\mathbf{D}}))$ is integrable over $\mathcal{D} \times [0,1]$, then Fubini's theorem implies that we can interchange the order of integration. Moreover, we have

$$\int_{\mathbf{D} \in \mathcal{D}} \int_{x \in [0,1]} [f(\mathbf{y}(\boldsymbol{\phi}_x, \mathbf{c}, \tilde{\mathbf{D}}))] dx \, d\mu(\mathbf{D}) = \int_{\mathbf{D} \in \mathcal{D}} \sum_{v=1}^{|\mathcal{V}|} f(\mathbf{y}(v, \mathbf{D})) \cdot |I_v(D)| d\mu(\mathbf{D})$$

$$= \int_{\mathbf{D} \in \mathcal{D}} \sum_{v=1}^{|\mathcal{V}|} f(\mathbf{y}(v, \mathbf{D})) \cdot \hat{\psi}_v(\mathbf{D}) d\mu(\mathbf{D}) = \langle \hat{\boldsymbol{\psi}}, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}}$$

Combing with (50), we have that

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\hat{\lambda}(\boldsymbol{\phi}) = \langle \hat{\boldsymbol{\psi}}, \boldsymbol{\eta}^{(0)} \rangle_{\mathcal{H}} \le \limsup_{k \to \infty} \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}))] d\lambda_k(\boldsymbol{\phi})$$

Similarly, for each $j \in \mathcal{N}$, we have that

$$\int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] d\hat{\lambda}(\boldsymbol{\phi}) = \int_{x \in [0,1]} \int_{\mathbf{D} \in \mathcal{D}} [R_j(s_j(\boldsymbol{\phi}_x, \mathbf{c}, \tilde{\mathbf{D}}), D_j)] d\mu(\mathbf{D}) dx$$

$$= \int_{\mathbf{D} \in \mathcal{D}} \int_{x \in [0,1]} [R_j(s_j(\boldsymbol{\phi}_x, \mathbf{c}, \tilde{\mathbf{D}}), D_j)] dx \, d\mu(\mathbf{D})$$

$$= \int_{\mathbf{D} \in \mathcal{D}} \sum_{v=1}^{|\mathcal{V}|} R_j(s_j(v, \mathbf{D}), D_j) \cdot |I_v(D)| d\mu(\mathbf{D}) = \langle \hat{\boldsymbol{\psi}}, \boldsymbol{\eta}^{(j)} \rangle_{\mathcal{H}}$$

$$\ge \liminf_{k \to \infty} \int_{\boldsymbol{\phi} \in \Phi} E_{\tilde{\mathbf{D}}}[R_j(s_j(\boldsymbol{\phi}, \mathbf{c}, \tilde{\mathbf{D}}), D_j)] d\lambda_k(\boldsymbol{\phi})$$

which completes our proof. $\quad\square$

*Proof of Claim 1:* We prove our result by contradiction. Suppose that $\mu(\hat{\mathcal{D}}) > 0$, then for each integer $p$, we define the set

$$\mathcal{D}_p = \{\mathbf{D} \in \mathcal{D} : \text{dist}(\hat{\boldsymbol{\psi}}(\mathbf{D}), \mathcal{L}) \ge \frac{1}{p}\}$$

where $\text{dist}(\hat{\boldsymbol{\psi}}(\mathbf{D}), \mathcal{L}) = \inf_{\mathbf{x} \in \mathcal{L}} \|\hat{\boldsymbol{\psi}}(\mathbf{D}) - \mathbf{x}\|_2^2$ denoting the distance from the point $\hat{\boldsymbol{\psi}}(\mathbf{D}) \in \mathbb{R}^{|\mathcal{V}|}$ to the set $\mathcal{L} \subset \mathbb{R}^{|\mathcal{V}|}$. Since the set $\mathcal{L}$ is closed, for any $\mathbf{D} \in \hat{\mathcal{D}} = \{\mathbf{D} \in \mathcal{D} : \hat{\boldsymbol{\psi}}(\mathbf{D}) \notin \mathcal{L}\}$, it holds that $\text{dist}(\hat{\boldsymbol{\psi}}(\mathbf{D}), \mathcal{L}) > 0$, which implies that

$$\hat{\mathcal{D}} = \bigcup_{p=1}^{\infty} \mathcal{D}_p$$

Moreover, note that $\mathcal{D}_1 \subset \mathcal{D}_2 \subset \cdots \subset \mathcal{D}_p \subset \ldots$, for each integer $p$, we define the set $\mathcal{E}_p = \mathcal{D}_p \setminus \mathcal{D}_{p-1} = \{\mathbf{D} \in \mathcal{D}_p : \mathbf{D} \notin \mathcal{D}_{p-1}\}$, where $\mathcal{D}_0 = \varnothing$. Obviously, the sets $\{\mathcal{E}_p\}$ are mutually disjoint and it holds that

$$\hat{\mathcal{D}} = \bigcup_{p=1}^{\infty} \mathcal{E}_p \quad \text{and} \quad \mathcal{D}_p = \bigcup_{l=1}^{p} \mathcal{E}_l \text{ for each integer } p$$

Then, from the coutable additivity of the measure $\mu$, we have

$$\mu(\hat{\mathcal{D}}) = \mu(\bigcup_{p=1}^{\infty} \mathcal{E}_p) = \sum_{p=1}^{\infty} \mu(\mathcal{E}_p) = \lim_{p \to \infty} \sum_{l=1}^{p} \mu(\mathcal{E}_l) = \lim_{p \to \infty} \mu(\bigcup_{l=1}^{p} \mathcal{E}_l) = \lim_{p \to \infty} \mu(\mathcal{D}_p)$$

Thus, we conclude that there exists an integer $p_0$, such that $\mu(\mathcal{D}_{p_0}) > 0$.

On the other hand, the sequence $\{\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}\}_{\forall k}$ converges strongly to $\hat{\boldsymbol{\psi}}$ implies that the sequence $\{\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}\}_{\forall k}$ converges to $\hat{\boldsymbol{\psi}}$ in measure, i.e., for any integer $p$,

$$\lim_{k \to \infty} \mu\left(\left\{\mathbf{D} \in \mathcal{D} : \|\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}(\mathbf{D}) - \hat{\boldsymbol{\psi}}(\mathbf{D})\|_2^2 \geq \frac{1}{p}\right\}\right) = 0$$

Moreover, note that $\mathcal{L}$ is a convex set, then for each $\mathbf{D} \in \mathcal{D}$, we must have $\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}(\mathbf{D}) \in \mathcal{L}$ for each $k$. Thus, it holds that

$$\text{dist}(\hat{\boldsymbol{\psi}}(\mathbf{D}), \mathcal{L}) \leq \|\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}(\mathbf{D}) - \hat{\boldsymbol{\psi}}(\mathbf{D})\|_2^2, \quad \forall k$$

which implies that the set $\mathcal{D}_p$ is contained in the set $\left\{\mathbf{D} \in \mathcal{D} : \|\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}(\mathbf{D}) - \hat{\boldsymbol{\psi}}(\mathbf{D})\|_2^2 \geq \frac{1}{p}\right\}$ for each $k$. As a result, we have that

$$\mu(\mathcal{D}_p) \leq \mu\left(\left\{\mathbf{D} \in \mathcal{D} : \|\frac{1}{k} \cdot \sum_{l=1}^{k} \boldsymbol{\psi}^{n_l}(\mathbf{D}) - \hat{\boldsymbol{\psi}}(\mathbf{D})\|_2^2 \geq \frac{1}{p}\right\}\right), \quad \forall k$$

and thus $\mu(\mathcal{D}_p) = 0$ for any integer $p$, which is a contradiction. $\square$

## Appendix D: Examples

We illustrate through the following examples that Corollary 1 can be used to find optimal capacity for problems studied in Mirchandani and Mishra (2002) and Swaminathan and Srinivasan (1999).

**Example 1.** Consider an inventory pooling example with one resource and two customers, i.e., $m = 1$ and $n = 2$. Assume the demand of the two customers are i.i.d. uniform distribution in $[0, 1]$. We are interested in Type I service levels with $\frac{\beta_1 + \beta_2}{2} \geq 75\%$ and the total inventory level $c \in [1, 2]$. For any fixed $\mathbf{w} \geq 0$, we assume without loss of generality that $w_1 \geq w_2$. Then solving the problem (16) gives us the optimal objective value as follows

$$\max_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j) = \begin{cases} w_1 + w_2 & \text{if } D_1 + D_2 \leq c \\ w_1 & \text{if } D_1 \leq c \text{ and } D_1 + D_2 > c \\ w_2 & \text{if } D_2 \leq c \text{ and } D_1 > c \\ 0 & \text{else} \end{cases}$$

Therefore, if $w_1 \geq w_2$, we have

$$E_{\tilde{\mathbf{D}}}\left[\max_{(\mathbf{y}, \mathbf{s}) \in P(\mathbf{c}, \mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j)\right]$$

$$= (w_1 + w_2) \int_0^1 \int_0^{\min\{c - D_1, 1\}} dD_2 dD_1 + w_1 \int_0^1 \int_{\min\{c - D_1, 1\}}^1 dD_2 dD_1 + w_2 \int_c^1 \int_0^c dD_2 dD_1$$

$$= w_1 + w_2(1 - \frac{1}{2}(2 - c)^2) \tag{52}$$

and

$$\max_{\mathbf{w} \geq 0} \sum_{j \in \mathcal{N}} w_j \beta_j - E_{\tilde{\mathbf{D}}}[\max_{(\mathbf{y},\mathbf{s}) \in P(\mathbf{c},\mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j)] = w_1 \beta_1 + w_2 \beta_2 - w_1 - w_2(1 - \frac{1}{2}(2-c)^2).$$

According to Corollary 1, $c$ is feasible if and only if

$$\max_{w_1 \geq w_2 \geq 0} \quad w_1(\beta_1 - 1) + w_2\left(\beta_2 - (1 - \frac{1}{2}(2-c)^2)\right) \leq 0.$$

Notice that $\beta_1 - 1 \leq 0$ and thus the maximum is attained when $w_1 = w_2$. Therefore, the condition above is equivalent to

$$\max_{w_2 \geq 0} \quad w_2\left(\beta_1 + \beta_2 - (2 - \frac{1}{2}(2-c)^2)\right) \leq 0,$$

or

$$2 - \frac{1}{2}(2-c)^2 \geq \beta_1 + \beta_2.$$

It follows immediately that in order to achieve individual Type I service levels $\beta_1$ and $\beta_2$ respectively, the minimum inventory is

$$c^* = 2 - 2\sqrt{1 - \frac{\beta_1 + \beta_2}{2}}.$$

This generalizes Theorem 2 of Swaminathan and Srinivasan (1999). For example, when $\beta_1 = \beta_2 = 80\%$, the minimum inventory level is 1.106.

We compare this with the minimum inventory level under a joint Type I service level $\beta$. For the same example, the minimum inventory $c$ should satisfy the constraint

$$\Pr\{D_1 + D_2 \leq c\} \geq \beta$$

which implies that

$$c \geq 2 - \sqrt{2 - 2\beta}.$$

When $\beta = 80\%$, the minimum inventory level is 1.368, which is 24% higher than that with individual 80% service constraints.

**Example 2.** Consider a W-system of the assemble-to-order model with two products and three components, i.e., $n = 2$ and $m = 3$. Components 1 and 2 are product-specific, while component 3 is common to both products. Under Type I service constraints, the problem (16) is specified as

$$\max_{(\mathbf{y},\mathbf{s}) \in P(\mathbf{c},\mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j) = \max \ w_1 z_1 + w_2 z_2 \tag{53}$$

$$\text{s.t.} \ z_1 D_1 + z_2 D_2 \leq c_3$$

$$z_j D_j \leq c_j \quad j = 1, 2,$$

$$z_j \in \{0, 1\} \quad j = 1, 2.$$

Similar to Example 1, we assume the demand of the two products are i.i.d. uniform distribution in $[0, 1]$. With this assumption, it is easy to see that $c_j \in [\beta_j, 1]$ for $j = 1, 2$, and $c_3 \leq c_1 + c_2$. We consider sufficiently high target service levels so that $c_3 \geq 1$. For any fixed $\mathbf{w} \geq 0$, we assume without loss of generality that $w_1 \geq w_2$. Then the optimal objective value of the problem (16) is

$$
\max_{(\mathbf{y},\mathbf{s}) \in P(\mathbf{c},\mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j) = \begin{cases} w_1 + w_2 & \text{if } D_1 + D_2 \leq c_3, \text{ and } D_1 \leq c_1, \text{ and } D_2 \leq c_2 \\ w_1 & \text{if } D_1 \leq c_1 \text{ and } D_2 > \min\{c_3 - D_1, c_2\} \\ w_2 & \text{if } D_2 \leq c_2 \text{ and } D_1 > c_1 \\ 0 & \text{else} \end{cases}
$$

Therefore, if $w_1 \geq w_2$, we have

$$
E_{\tilde{\mathbf{D}}}[\max_{(\mathbf{y},\mathbf{s}) \in P(\mathbf{c},\mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j)]
$$

$$
= (w_1 + w_2) \int_0^{c_1} \int_0^{\min\{c_3 - D_1, c_2\}} dD_2 dD_1 + w_1 \int_0^{c_1} \int_{\min\{c_3 - D_1, c_2\}}^1 dD_2 dD_1 + w_2 \int_{c_1}^1 \int_0^{c_2} dD_2 dD_1
$$

$$
= w_1 c_1 + w_2 (c_2 - \frac{1}{2}(c_1 + c_2 - c_3)^2) \tag{54}
$$

and

$$
\max_{\mathbf{w} \geq 0} \sum_{j \in \mathcal{N}} w_j \beta_j - E_{\tilde{\mathbf{D}}}[\max_{(\mathbf{y},\mathbf{s}) \in P(\mathbf{c},\mathbf{D})} \sum_{j \in \mathcal{N}} w_j R_j(s_j, D_j)] = w_1 \beta_1 + w_2 \beta_2 - w_1 c_1 - w_2 (c_2 - \frac{1}{2}(c_1 + c_2 - c_3)^2).
$$

According to Corollary 1, $c$ is feasible if and only if

$$
\max_{w_1 \geq w_2 \geq 0} \quad w_1(\beta_1 - c_1) + w_2 \left( \beta_2 - (c_2 - \frac{1}{2}(c_1 + c_2 - c_3)^2) \right) \leq 0.
$$

Since $\beta_1 - c_1 \leq 0$, the condition is equivalent to

$$
\max_{w_2 \geq 0} \quad w_2 \left( \beta_1 + \beta_2 - (c_1 + c_2 - \frac{1}{2}(c_1 + c_2 - c_3)^2) \right) \leq 0,
$$

or

$$
c_1 + c_2 - \frac{1}{2}(c_1 + c_2 - c_3)^2 \geq \beta_1 + \beta_2.
$$

Notice that the condition depends on $\bar{c} = c_1 + c_2$, but not the individual values of $c_1$ and $c_2$. Therefore, we can derive the minimum inventory level by solving

$$
\min \quad \bar{c} + c_3
$$

$$
\text{s.t. } \bar{c} - \frac{1}{2}(\bar{c} - c_3)^2 \geq \beta_1 + \beta_2
$$

$$
\bar{c} \leq 2
$$

$$
\bar{c} \geq c_3.
$$

The first two constraints implies that

$$c_3 \geq 2 - 2\sqrt{1 - \frac{\beta_1 + \beta_2}{2}}.$$

Comparing this with the optimal inventory level in Example 1, we see that everything else being equal, the common component in the W system holds more inventory than that in single-resource pooling. Solving the optimization problem, we derive that the minimum total inventory level of the three components is

$$\bar{c} + c_3 = \begin{cases} 2(\beta_1 + \beta_2) - \frac{1}{4} & \text{if } \beta_1 + \beta_2 \leq \frac{15}{8} \\ 4 - \sqrt{4 - 2(\beta_1 + \beta_2)} & \text{if } \beta_1 + \beta_2 > \frac{15}{8} \end{cases}$$

This is a special case of the main result of Mirchandani and Mishra (2002). A straightforward extension of this analysis would give a different proof of Theorem 2 of Mirchandani and Mishra (2002).

### Appendix E: Proof of Theorem 3

We first present the two well-known results, which will be useful for our proof of Theorem 3.

LEMMA 4. *(Azuma's Inequality (Azuma, 1967)) Suppose $\{X_k, k = 0, 1, 2, \ldots\}$ is a martingale and $|X_k - X_{k-1}| < c_k$ almost surely for each $k$, then for all positive integer $N$ and all positive real $\epsilon$,*

$$P(|X_N - X_0| > \epsilon) \leq 2 \cdot \exp(\frac{-\epsilon^2}{2 \sum_{k=1}^{N} c_k^2})$$

LEMMA 5. *(Borel-Cantelli Lemma (Borel, 1909)) Let $E_1, E_2, \ldots$ be a sequence of events in a probability space. If*

$$\sum_{n=1}^{\infty} P(E_n) < \infty$$

*then*

$$P\left(\bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} E_k\right) = 0$$

*where $\bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} E_k$ denotes the set of outcomes that occur infinite times with the sequence $\{E_k\}_{k \geq 1}$.*

Now we are ready to prove Theorem 3. The main idea of the proof is to show that by following the dual update step of the Max-Weighted-Service policy, the gap between the expected cost of our policy and the optimal value of (11), as well as the gap between the achieved and the target service level, can both be bounded by some functions of the dual variables, which diminish under carefully chosen step sizes, as $T \to \infty$.

*Proof of Theorem 3:* We first prove (27). From weak duality (Shapiro, 2001), it holds that $G(\mathbf{w}^*) \leq \text{Obj}$ (11). Thus, it is enough to prove that

$$\limsup_{T \to \infty} \frac{1}{T} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t}[f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t))] \leq G(\mathbf{w}^*)$$

holds almost surely. Note that by definition of $G(\cdot)$ and $\mathbf{w}^*$, we have that

$$G(\mathbf{w}^*) = \max_{\mathbf{w} \geq 0} G(\mathbf{w}) \geq \frac{1}{T} \cdot \sum_{t=1}^{T} G(\mathbf{w}^{(t)})$$

$$\geq \frac{1}{T} \cdot \sum_{t=1}^{T} \left( \sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \beta_j + \mathrm{E}_{\tilde{\mathbf{D}}^t}[f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t))] - \sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t}[R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t)] \right)$$

Re-arranging terms, we have

$$\frac{1}{T} \cdot \sum_{t=1}^{T} \mathrm{E}_{\tilde{\mathbf{D}}^t}[f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t))] - G(\mathbf{w}^*) \leq \frac{1}{T} \cdot \sum_{t=1}^{T} \left( -\sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \beta_j + \sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t}[R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t)] \right) \tag{55}$$

We now proceed to upper bound the right hand side of the above inequality. From the update rule (26), we have that

$$\|\mathbf{w}^{(t+1)}\|^2 \leq \left\| \left( w_j^{(t)} + \gamma_T \cdot \left( \beta_j - R_j \left( s_j \left( \boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)} \right), D_j^{(t)} \right) \right), \ j \in \mathcal{N} \right) \right\|^2$$

$$= \|\mathbf{w}^{(t)}\|^2 + \gamma_T^2 \cdot \left\| \left( \beta_j - R_j \left( s_j \left( \boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)} \right), D_j^{(t)} \right), \ j \in \mathcal{N} \right) \right\|^2$$

$$+ 2\gamma_T \cdot \sum_{j \in \mathcal{N}} \left( w_j^{(t)} \cdot \mathrm{E}_{\mathbf{D}^{(t)}} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right] - w_j^{(t)} \cdot R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right)$$

$$+ 2\gamma_T \cdot \sum_{j \in \mathcal{N}} \left( w_j^{(t)} \cdot \beta_j - w_j^{(t)} \cdot \mathrm{E}_{\mathbf{D}^{(t)}} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right] \right) \tag{56}$$

where $\|\cdot\|$ denotes the $L_2$ norm and we denote $(a_j, \ j \in \mathcal{N})$ as a $n$ dimensional vector with $a_j$ on its $j$-th component for any $j \in \mathcal{N}$. Moreover, for each $t$, we denote

$$L_t = \frac{2}{T\gamma_T} \cdot \sum_{j \in \mathcal{N}} \left( w_j^{(t)} \cdot \mathrm{E}_{\mathbf{D}^{(t)}} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right] - w_j^{(t)} \cdot R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right) \tag{57}$$

Obviously, $\{L_t\}_{t \geq 1}$ is a sequence of martingale difference with respect to the filtration $\{\mathcal{F}_t\}_{t \geq 1}$, where for each $t$, $\mathcal{F}_t$ denotes the $\sigma$-algebra $\sigma\left(\mathbf{D}^{(1)}, \ldots, \mathbf{D}^{(t)}\right)$. From the update rule (26), we have $w_j^{(t+1)} \leq t \cdot \gamma_T \cdot \beta_j$ for each $j \in \mathcal{N}$ and each $t$. Thus, by Assumption 3, there exists a constant $\hat{C}_1$ such that $|L_t| \leq \hat{C}_1$ almost surely for each $t$. Note that it follows from Assumption 3 that for any $\mathbf{w}^{(t)}$ and $\tilde{\mathbf{D}}$, we have

$$\left\| \left( \beta_j - R_j \left( s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j \right), \ j \in \mathcal{N} \right) \right\|_2^2 \leq n(1+C)^2. \tag{58}$$

Further note that for each $t$, we have

$$w_j^{(t)} \cdot \mathrm{E}_{\mathbf{D}^{(t)}} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right] = w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t) \right] \quad (\text{since } \tilde{\mathbf{D}}^t \sim \mathbf{D}^{(t)})$$

Then, (56) implies that

$$\|\mathbf{w}^{(t+1)}\|^2 \leq \|\mathbf{w}^{(t)}\|^2 + \gamma_T^2 \cdot n(1+C)^2 + T\gamma_T^2 \cdot L_t + 2\gamma_T \cdot \sum_{j \in \mathcal{N}} \left( w_j^{(t)} \cdot \beta_j - w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t) \right] \right)$$

and thus, by re-arranging terms, we have

$$\sum_{j \in \mathcal{N}} \left( -w_j^{(t)} \cdot \beta_j + w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t) \right] \right) \geq \frac{1}{2\gamma_T} \left( -\|\mathbf{w}^{(t+1)}\|^2 + \|\mathbf{w}^{(t)}\|^2 + \gamma_T^2 \cdot n(1+C)^2 + T\gamma_T^2 \cdot L_t \right) \tag{59}$$

Plugging (59) into (55), we have that

$$
\begin{aligned}
&\frac{1}{T} \cdot \sum_{t=1}^{T} \mathrm{E}_{\tilde{\mathbf{D}}^t}[f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t))] - G(\mathbf{w}^*) \\
&\leq \frac{1}{T} \cdot \sum_{t=1}^{T} \left( -\sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \beta_j + \sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}^t}[R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t)] \right) \\
&\leq \frac{1}{2T\gamma_T} \cdot \sum_{t=1}^{T} \left\{ \|\mathbf{w}^{(t)}\|^2 - \|\mathbf{w}^{(t+1)}\|^2 + \gamma_T^2 \cdot n(1+C)^2 + T\gamma_T^2 \cdot L_t \right\} \\
&= \frac{1}{2T\gamma_T} \cdot (\|\mathbf{w}^{(1)}\|^2 - \|\mathbf{w}^{(T+1)}\|^2) + \frac{\gamma_T \cdot n(1+C)^2}{2} + \frac{\gamma_T}{2} \cdot \sum_{t=1}^{T} L_t \\
&\leq \frac{\gamma_T \cdot n(1+C)^2}{2} + \frac{\gamma_T}{2} \cdot \sum_{t=1}^{T} L_t
\end{aligned}
\tag{60}
$$

where the last inequality holds since $\mathbf{w}^{(1)} = 0$.

For any $a > 0$, define $E_T(a)$ as the event that $\gamma_T \cdot |\sum_{t=1}^{T} L_t| \geq a$. Since we have shown $|L_t| \leq \hat{C}_1$ almost surely for each $t$ and $\gamma_T = T^{-(\frac{1}{2}+\epsilon)}$ for some $\epsilon \in (0, 1/2)$, by Azuma's inequality, we have

$$P(E_T(a)) = P(\gamma_T \cdot |\sum_{t=1}^{T} L_t| \geq a) \leq 2 \cdot \exp(-\frac{a^2}{2\hat{C}_1^2} \cdot T^{2\epsilon})$$

Note that the above inequality implies that

$$\sum_{T=1}^{\infty} P(E_T(a)) \leq 2 \cdot \sum_{T=1}^{\infty} \exp(-\frac{a^2}{2\hat{C}_1^2} \cdot T^{2\epsilon}) < \infty$$

Then, by Borel-Cantelli Lemma, we know that $P\left(\bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} E_k(a)\right) = 0$ for each $a > 0$. Thus, it holds that

$$\lim_{T \to \infty} \gamma_T \cdot \sum_{t=1}^{T} L_t = 0 \quad \text{almost surely}$$

when $\gamma_T = T^{-(\frac{1}{2}+\epsilon)}$ for some $\epsilon \in (0, 1/2)$, which completes our proof of (27).

We then prove (28). To that end, we define, for each $j \in \mathcal{N}$ and each $\tau = 1, \ldots, T$,

$$\rho_{j,\tau} = \beta_j - \frac{1}{\tau} \sum_{t=1}^{\tau} R_j(s_j(\phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}^t), \tilde{D}_j^t).$$

Then (28) follows if $\limsup_{T \to \infty} \rho_{j,T} \leq 0$ almost surely for all $j \in \mathcal{N}$. To that end, we first prove that $\tau \rho_{j,\tau} \preceq \frac{w_j^{(\tau+1)}}{\gamma_T}$, $\forall j \in \mathcal{N}$ for any $\tau \geq 1$, where $a \preceq b$ denotes random variable $a$ is first-order stochastic dominated by random variable $b$. Then, in Lemma 6 below, we show that for each $j \in \mathcal{N}$, $\limsup_{T \to \infty} \frac{w_j^{(T+1)}}{T \gamma_T} = 0$ almost surely.

We prove $\tau \rho_{j,\tau} \preceq \frac{1}{\gamma_T} \cdot w_j^{(\tau+1)}$ by induction. When $\tau = 1$, noticing that $\mathbf{w}^{(1)} = 0$, we have that

$$\begin{aligned}
\rho_{j,1} &= \beta_j - R_j \left( s_j(\phi_{\mathbf{w}^{(1)}}, \mathbf{c}, \tilde{\mathbf{D}}^1), \tilde{D}_j^1 \right) \\
&= w_j^{(1)} + \beta_j - \left[ R_j \left( s_j(\phi_{\mathbf{w}^{(1)}}, \mathbf{c}, \tilde{\mathbf{D}}^1), \tilde{D}_j^1 \right) \right] \\
&\preceq \left[ w_j^{(1)} + \beta_j - R_j \left( s_j(\phi_{\mathbf{w}^{(1)}}, \mathbf{c}, \mathbf{D}^{(1)}), D_j^{(1)} \right) \right]^+ \quad \text{(since } \mathbf{D}^{(1)} \sim \tilde{\mathbf{D}}^1 \text{)} \\
&= \frac{1}{\gamma_T} \cdot w_j^{(2)}
\end{aligned}$$

Now assume that we have $(\tau - 1) \rho_{j,\tau-1} \preceq \frac{1}{\gamma_T} \cdot w_j^{(\tau)}$, then

$$\begin{aligned}
\tau \rho_{j,\tau} &= (\tau - 1) \rho_{j,\tau-1} + \beta_j - R_j \left( s_j(\phi_{\mathbf{w}^{(\tau)}}, \mathbf{c}, \tilde{\mathbf{D}}^\tau), \tilde{D}_j^\tau \right) \\
&\preceq (\tau - 1) \rho_{j,\tau-1} + \beta_j - R_j \left( s_j(\phi_{\mathbf{w}^{(\tau)}}, \mathbf{c}, \mathbf{D}^{(\tau)}), \tilde{D}_j^{(\tau)} \right) \quad \text{(since } \mathbf{D}^{(\tau)} \sim \tilde{\mathbf{D}}^\tau \text{)} \\
&\preceq \frac{1}{\gamma_T} \cdot w_j^{(\tau)} + \beta_j - R_j \left( s_j(\phi_{\mathbf{w}^{(\tau)}}, \mathbf{c}, \mathbf{D}^{(\tau)}), \tilde{D}_j^{(\tau)} \right) \\
&\leq \left[ \frac{1}{\gamma_T} \cdot w_j^{(\tau)} + \beta_j - R_j \left( s_j(\phi_{\mathbf{w}^{(\tau)}}, \mathbf{c}, \mathbf{D}^{(\tau)}), \tilde{D}_j^{(\tau)} \right) \right]^+ \\
&= \frac{1}{\gamma_T} \cdot w_j^{(\tau+1)}
\end{aligned}$$

Thus $T \rho_{j,T} \preceq \frac{1}{\gamma_T} \cdot w_j^{(T+1)}$, $\forall j \in \mathcal{N}$, which completes the proof. $\qquad \square$

LEMMA 6. *If $\mathbf{c}$ is feasible, then for any $j \in \mathcal{N}$, $\limsup_{T \to \infty} \frac{w_j^{(T+1)}}{T \gamma_T} = 0$ almost surely.*

*Proof:* It follows from (56) and Assumption 3 that

$$\begin{aligned}
\|\mathbf{w}^{(t+1)}\|^2 \leq{}& \|\mathbf{w}^{(t)}\|^2 + \gamma_T^2 \cdot n(1+C)^2 + T \gamma_T^2 \cdot L_t \\
&+ 2\gamma_T \cdot \sum_{j \in \mathcal{N}} \left( w_j^{(t)} \cdot \beta_j - w_j^{(t)} \cdot \mathrm{E}_{\mathbf{D}^{(t)}} \left[ R(s_j(\phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right] \right)
\end{aligned}$$

Notice that

$$\mathrm{E}_{\mathbf{D}^{(t)}} \left[ \sum_{j \in \mathcal{N}} w_j^{(t)} \cdot R_j(s_j(\phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \mathbf{D}^{(t)}), D_j^{(t)}) \right] = \sum_{j \in \mathcal{N}} w_j^{(t)} \cdot \mathrm{E}_{\tilde{\mathbf{D}}} \left[ R_j(s_j(\phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j) \right] \quad (\mathbf{D}^{(t)} \sim \tilde{\mathbf{D}})$$

Also, from weak duality, it holds $G(\mathbf{w}^{(t)}) \leq G(\mathbf{w}^*) \leq \text{Obj (11)}$. Then, we have that

$$\sum_{j \in \mathcal{N}} w_j^{(t)} \beta_j + \mathrm{E}_{\tilde{\mathbf{D}}} \left[ f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}})) \right] - \sum_{j \in \mathcal{N}} w_j^{(t)} \, \mathrm{E}_{\tilde{\mathbf{D}}} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j) \right] \leq \text{Obj (11)}$$

When the capacity level $\mathbf{c}$ is feasible, the objective value of (11) is finite and we denote $\hat{C}_2$ as its upper bound. Thus, we have that

$$\sum_{j \in \mathcal{N}} w_j^{(t)} \beta_j - \sum_{j \in \mathcal{N}} w_j^{(t)} \, \mathrm{E}_{\tilde{\mathbf{D}}} \left[ R_j(s_j(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j) \right] \leq \text{Obj (11)} - \mathrm{E}_{\tilde{\mathbf{D}}} \left[ f(\mathbf{y}(\boldsymbol{\phi}_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}})) \right] \leq \hat{C}_2 \quad (61)$$

Therefore, we must have

$$\|\mathbf{w}^{(t+1)}\|_2^2 \leq \|\mathbf{w}^{(t)}\|_2^2 + \gamma_T^2 \cdot n(1+C)^2 + 2\gamma_T \cdot \hat{C}_2 + T\gamma_T^2 \cdot L_t \quad \forall t = 1, 2, \ldots, T \qquad (62)$$

Summing inequality (62) from $t = 1$ to $T$, we get

$$\|\mathbf{w}^{(T+1)}\|_2^2 \leq \|\mathbf{w}^{(1)}\|_2^2 + T\gamma_T^2 \cdot n(1+C)^2 + 2T\gamma_T \cdot \hat{C} + T\gamma_T^2 \cdot \sum_{t=1}^{T} L_t$$

Then, we have

$$\frac{1}{T^2 \gamma_T^2} \cdot \|\mathbf{w}^{(T+1)}\|_2^2 \leq \frac{1}{T} \cdot n(1+C)^2 + \frac{2\hat{C}}{T\gamma_T} + \frac{1}{T} \cdot \sum_{t=1}^{T} L_t \qquad (63)$$

Since we have shown $|L_t| \leq \hat{C}_1$ almost surely for each $t$, we can again apply the combination of Azuma's inequality and the Borel-Cantelli lemma to show that

$$\lim_{T \to \infty} \frac{1}{T} \cdot \sum_{t=1}^{T} L_t = 0 \quad \text{almost surely}$$

Thus, we have for each $j \in \mathcal{N}$,

$$\limsup_{T \to \infty} \frac{w_j^{(T+1)}}{T\gamma_T} = 0 \quad \text{almost surely}$$

which completes the proof. $\quad \square$

*Proof of Corollary 2.* The proof of Corollary 2 follows the same main steps as that of Theorem 3 with minor modifications outlined below. Notice that Corollary 2, which focuses only on the expected performance of Algorithm 1, is weaker than Theorem 3. Thus, Assumption 3, which is needed for the proof of Theorem 3, is now replaced with a weaker condition, i.e., $E_D[R_j(s_j, \tilde{D}_j)^2] \leq C$ for all $s_j \geq 0$.

In particular, Assumption 3 is only used in Theorem 3 to prove the boundedness of $L_t$ and inequality (58). For the proof of Corollary 2, we can take expectation over $w^{(t)}$ and $D^{(t)}$ on both sides of (56) and (58). Then (56) implies that $E_{w^{(t)}, D^{(t)}}[L_t] = 0$ and thus bounded for each $t$ . Also,

inequality (58) holds in expectation as long as $E_D[R_j(s_j, \tilde{D}_j)^2] \leq C$ for all $s_j \geq 0$, which is the condition assumed in Corollary 2.

Then, (60) together with $\mathrm{E}_{\mathbf{w}^{(t)}, \tilde{\mathbf{D}}^{(t)}}[L_t] = 0$, directly imply the convergence rate on the expected allocation cost should be $O(\gamma_T)$.

Finally, (63) implies, for any $j$,

$$\mathrm{E}[\frac{w_j^{(T+1)}}{T\gamma_T}] \leq O(\max\{\sqrt{\frac{1}{T}}, \sqrt{\frac{1}{T\gamma_T}}\}).$$

Thus, we have

$$\mathrm{E}[\rho_{j,T}] \leq \mathrm{E}[\frac{w_j^{(T+1)}}{T\gamma_T}] \leq O(\max\{\sqrt{\frac{1}{T}}, \sqrt{\frac{1}{T\gamma_T}}\}),$$

which proves the convergence rate on the expected service level. $\qquad\square$

## Appendix F: Discussion over the polymatroid assumption

A wide range of capacity allocation problems enjoys the polymatorid structure as illustrated below.

- In single-resource pooling, $Q(c, \mathbf{D})$ is represented by (1). It is straightforward to see that $Q(c, \mathbf{D})$ can be reformulated as (29) with the corresponding submodular set function

$$q(U|c, \mathbf{D}) = \min\{c, \sum_{j \in U} D_j\}.$$

- Consider a capacity planning problem of a more general directed network $G = (V, E)$ where $V$ is the set of nodes and $E$ is the set of arcs with $|E| = m$. There is a single supply node $u \in V$ with unlimited capacity and a set of demand nodes $\mathcal{N} \subset V$ whose demands are random. The problem is to decide the capacity $c_e$ of each arc $e \in E$ in anticipation of random demands, and route the supply from $u$ through the arcs in $E$ to satisfy the demand of $\mathcal{N}$ after demand realization. For this problem, the feasible set $Q(\mathbf{c}, \mathbf{D})$ is

$$\{\mathbf{s} \in \mathbb{R}_+^{\mathcal{N}} : \sum_{j \in U} s_j \leq q(U|\mathbf{c}, \mathbf{D}), \quad \forall U \subseteq \mathcal{N}\}$$

  where $q(U|\mathbf{c}, \mathbf{D})$ is the total maximum flow from $u$ to the demand nodes in $U$ given the capacity of arcs $\mathbf{c} = (c_e : e \in E)$ and the demand realization $\mathbf{D} = (D_j : j \in \mathcal{N})$. It is well known that $Q(\mathbf{c}, \mathbf{D})$ is a polymatroid (Megiddo, 1974; He et al., 2012). Specifically, when the network is bipartite, the problem reduces to the process flexibility problem.

- In a special case of assemble-to-order system where $Q(\mathbf{c}, \mathbf{D})$ is represented by (3) and the component consumption matrix is represented by (4), the corresponding submodular function is

$$q(U|\mathbf{c}, \mathbf{D}) = \min\{c_{n+1}, \sum_{j \in U} \min\{c_j, D_j\}\}.$$

In all these examples, the submodular function $q(U|\mathbf{c}, \mathbf{D})$ for any $U \subseteq \mathcal{N}$ can be defined as the maximum total fulfilled demand of all customers in the subset $U$ using all available capacity $\mathbf{c}$. Given the above notations, we prove Theorem 4.

*Proof of Theorem 4:* Given the randomly selected weight vector $\mathbf{w} = \tilde{\mathbf{w}}$, and the realized demand $\mathbf{D}$, the Max-Weighted-Service policy fulfills the demands by solving the following Max-Weighted-Service problem (16)

$$\max \sum_{j \in \mathcal{N}} w_j \cdot a_j(D_j)s_j + w_j b_j(D_j) - v_j s_j \tag{64}$$
$$\text{s.t. } \mathbf{s} \in Q(\mathbf{c}, \mathbf{D})$$

By Assumption 4, the feasible set $Q(\mathbf{c}, \mathbf{D})$ is a polymatroid. We denote the submodular set function that defines the polymatroid $Q(\mathbf{c}, \mathbf{D})$ by $q(U|\mathbf{c}, \mathbf{D})$ and assume that

$$w_1 \cdot a_1(D_1) - v_1 \geq w_2 \cdot a_2(D_2) - v_2 \geq \cdots \geq w_n \cdot a_n(D_n) - v_n.$$

Then it is well-known (Welsh, 2010) that the following solution is optimal to problem (64):

$$s_1^* = q(\{1\}|\mathbf{c}, \mathbf{D})$$
$$s_j^* = q(\{1, 2, \cdots, j\}|\mathbf{c}, \mathbf{D}) - q(\{1, 2, \cdots, j-1\}|\mathbf{c}, \mathbf{D}), \quad j = 2, \cdots, n$$

That is, the customers are fulfilled according to a non-increasing order of $w_j \cdot a_j(D_j) - v_j$, which is an index policy. $\square$

We now show that when the service measure function $R_j$ represents Type II service levels for each $j \in \mathcal{N}$, a randomized anticipative index policy is not just asymptotically optimal, but actually optimal. Indeed, for each $T$, we denote $\tilde{\lambda}_T$ as the uniform distribution over $\{\phi_{\mathbf{w}^{(1)}}, \ldots, \phi_{\mathbf{w}^{(T)}}\}$, where $\{\mathbf{w}^{(1)}, \ldots, \mathbf{w}^{(T)}\}$ denotes the sequence of weights generated by Algorithm 1. From Theorem 4, we know that each $\phi_{\mathbf{w}^{(t)}}$ can be characterized as a deterministic anticipative index policy. Obviously, the total number of index lists is finite, denoted as $\hat{L}$. Then, for each $T$, $\tilde{\lambda}_T$ can be regarded as a $\hat{L}$-dimensional vector in a compact set, thus the sequence $\{\tilde{\lambda}_T\}_{T \to \infty}$ must have a convergent subsequence, and we denote $\tilde{\lambda}$ as the limit of this subsequence. Clearly, $\hat{\lambda}$ denotes a randomized anticipative index policy. Then from Theorem 3, the following two inequalities hold almost surely:

$$\int_{\phi \in \Phi} \mathrm{E}_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi, \mathbf{c}, \tilde{\mathbf{D}}))]d\tilde{\lambda}(\phi) \leq \limsup_{T \to \infty} \frac{1}{T} \cdot \sum_{t=1}^{T} \mathrm{E}_{\tilde{\mathbf{D}}}[f(\mathbf{y}(\phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}))] \leq \text{Obj (11)}$$

and

$$\int_{\phi \in \Phi} \mathrm{E}_{\tilde{\mathbf{D}}}[R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}))]d\tilde{\lambda}(\phi) \geq \liminf_{T \to \infty} \frac{1}{T} \cdot \sum_{t=1}^{T} \mathrm{E}_{\tilde{\mathbf{D}}}\left[R_j(s_j(\phi_{\mathbf{w}^{(t)}}, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)\right] \geq \beta_j, \quad \forall j \in \mathcal{N}.$$

Thus, we conclude that $\tilde{\lambda}$ is an optimal solution to (11) almost surely.

Although $\tilde{\lambda}$ is optimal to the single period formulation (11) with probability 1, in practice, finding the optimal solution $\tilde{\lambda}$ may require solving a large-scale linear programming and thus it may be computationally inefficient. Instead, Theorem 3 shows that one can use the uniform distribution $\tilde{\lambda}_T$ to approximate $\tilde{\lambda}$, which is asymptotically optimal as $T \to \infty$.

## Appendix G: Proof of Theorem 5

*Proof of Theorem 5:* We denote by $\mathrm{Obj}_{\mathbf{c}}(11)$ the objective value of (11) given capacity level $\mathbf{c}$. From Theorem 2, we have that

$$\max_{\mathbf{w} \geq 0} H(\mathbf{w}, \mathbf{c}) = \begin{cases} p(\mathbf{c}) + \mathrm{Obj}_{\mathbf{c}}(11), & \text{if (11) is feasible for } \mathbf{c} \\ +\infty, & \text{if (11) is infeasible for } \mathbf{c} \end{cases}$$

Denote by $\mathbf{c}^*$ one optimal solution of (10). Then, we have

$$\mathrm{Obj}\ (10) = p(\mathbf{c}^*) + \mathrm{Obj}_{\mathbf{c}^*}(11) = \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \mathbf{c}^*) \geq \min_{\mathbf{c} \geq 0} \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \mathbf{c}) \tag{65}$$

Denote by $\hat{\mathbf{c}}$ one optimal solution of (31). Then, we have that

$$\min_{\mathbf{c} \geq 0} \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \mathbf{c}) = \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \hat{\mathbf{c}}) < +\infty$$

which implies that (11) is feasible under the capacity level $\hat{c}$. Thus, we have

$$\min_{\mathbf{c} \geq 0} \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \mathbf{c}) = \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \hat{\mathbf{c}}) = p(\hat{\mathbf{c}}) + \mathrm{Obj}_{\hat{\mathbf{c}}}(11) \geq \mathrm{Obj}\ (10) \tag{66}$$

As a result, we have $\mathrm{Obj}\ (10) = \min_{\mathbf{c} \geq 0} \max_{\mathbf{w} \geq 0}\ H(\mathbf{w}, \mathbf{c})$, and all the inequalities in (65) and (66) hold as equality, which implies that $\mathbf{c}^*$ is an optimal solution to (31) and $\hat{\mathbf{c}}$ is an optimal solution to (10). $\quad\square$

## Appendix H: Proof of Lemma 2

*Proof of Lemma 2.* By definition, for any $\mathbf{w} \geq 0$ and $\mathbf{D}$,

$$g(\mathbf{w}, \mathbf{c}; \mathbf{D}) = \min\ f(\mathbf{y}) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(s_j, D_j) \tag{67}$$
$$\text{s.t.}\ \ (\mathbf{s}, \mathbf{y}) \in P(\mathbf{c}, \mathbf{D})$$

Let $(s^*(\mathbf{c}), y^*(\mathbf{c}))$ be an optimal solution. (Here we assume that $\mathbf{w}$ and $\mathbf{D}$ are fixed and thus drop the dependence on them in $(s^*(\mathbf{c}), y^*(\mathbf{c}))$.) Then we have $g(\mathbf{w}, \mathbf{c}; \mathbf{D}) = f(\mathbf{y}) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(s_j^*(\mathbf{c}), \mathbf{D})$. For any two points $\mathbf{c}^1, \mathbf{c}^2$ and any constant $0 < \alpha < 1$,

$$(\alpha s^*(\mathbf{c}^1) + (1-\alpha)s^*(\mathbf{c}^2)), ((\alpha y^*(\mathbf{c}^1) + (1-\alpha)y^*\mathbf{c}^2)))$$

must be a feasible solution to (67) when $\mathbf{c} = \alpha\mathbf{c}^1 + (1 - \alpha)\mathbf{c}^2$. Then we have $(\alpha s^*(\mathbf{c}^1) + (1 - \alpha)s^*(\mathbf{c}^2), \alpha y^*(\mathbf{c}^1) + (1 - \alpha)y^*(\mathbf{c}^2)) \in P((\alpha\mathbf{c}^1 + (1 - \alpha)\mathbf{c}^2), \mathbf{D})$. Thus, from the concavity of $R_j(s_j, D_j)$ in $s_j$, we have that

$$\alpha g(\mathbf{w}, \mathbf{c}^1; \mathbf{D}) + (1 - \alpha)g(\mathbf{w}, \mathbf{c}^2; \mathbf{D})$$

$$= (\alpha \cdot f(\mathbf{y}^*(\mathbf{c}^1)) + (1 - \alpha) \cdot f(\mathbf{y}^*(\mathbf{c}^2))) - \sum_{j \in \mathcal{N}} w_j \cdot (\alpha R_j(s_j^*(\mathbf{c}^1), D_j) + (1 - \alpha)R_j(s_j^*(\mathbf{c}^2), D_j))$$

$$\geq f(\alpha \cdot \mathbf{y}^*(\mathbf{c}^1) + (1 - \alpha) \cdot \mathbf{y}^*(\mathbf{c}^2)) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(\alpha s_j^*(\mathbf{c}^1) + (1 - \alpha)s_j^*(\mathbf{c}^2), D_j)$$

$$\geq \min_{(\mathbf{s}, \mathbf{y}) \in P(\alpha\mathbf{c}^1 + (1-\alpha)\mathbf{c}^2, \mathbf{D})} f(\mathbf{y}) - \sum_{j \in \mathcal{N}} w_j \cdot R_j(s_j, D_j)$$

$$= g(\mathbf{w}, \alpha\mathbf{c}^1 + (1 - \alpha)\mathbf{c}^2; \mathbf{D})$$

We conclude that $g(\mathbf{w}, \mathbf{c}; \tilde{\mathbf{D}})$ is convex in $\mathbf{c}$ for any $\mathbf{w} \geq 0$ and $\tilde{\mathbf{D}}$, and thus $H(\mathbf{w}, \mathbf{c})$ is a convex function of $\mathbf{c}$ for any $\mathbf{w} \geq 0$.    $\square$

### Appendix I: Asymptotic Bound

In this section, we consider inventory pooling with i.i.d. demand distribution, equivalent target service level $\beta$, and Type I service constraints. We show that Corollary 1 can be used to obtain a closed-form expression of the asymptotically optimal capacity level per customer when the number of customers goes to infinity.

THEOREM 7. *For inventory pooling with $n$ customers, assume that the demands are i.i.d. with a common strictly increasing continuous distribution function $F(\cdot)$ and the Type I service level targets are all equal to $\beta$. Then the per customer capacity level $c^*$ is asymptotic optimal as $n \to \infty$, where*

$$c^* = \beta \cdot \max_{\xi}\{\xi - \frac{1}{\beta} \ E_{\tilde{D}}[(\xi - \tilde{D})^+]\}.$$

*Morover, if the distribution function $F(\cdot)$ has a finite mean $\mu$ and a finite standard deviation $\sigma$, then it holds that*

$$\beta\mu - \sqrt{\beta(1 - \beta)}\sigma \leq c^* \leq \beta\mu \tag{68}$$

*Proof:* From Corollary 1, a given capacity level $c$ is feasible is equivalent to the following condition:

$$\max_{w \geq 0} \sum_j w_j\beta_j - \ E_{\tilde{\mathbf{D}}}[\max_{\phi \in \Phi} \sum_j w_j R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \leq 0 \tag{69}$$

From the symmetry of the problem, (69) must have an optimal solution $(w_j^* : j = 1, ..., n)$ such that $w_j^* = w_1^*$ for any $j \neq 1$. Thus, condition (69) is equivalent to

$$n \cdot \beta - \ E_{\tilde{\mathbf{D}}}[\max_{\phi \in \Phi} \sum_j R_j(s_j(\phi, \mathbf{c}, \tilde{\mathbf{D}}), \tilde{D}_j)] \leq 0 \tag{70}$$

Note that by definition of $R_j$ for Type I service constraints, the inner maximization problem in (70) can be expressed as

$$M_n^*(c) := \max\{|A| : \sum_{j \in A \subset \mathcal{N}} \tilde{D}_j \le c\}$$

and condition (70) is equivalent to

$$\beta \le \frac{1}{n} \cdot \mathrm{E}_{\tilde{\mathbf{D}}}[M_n^*(c)] \tag{71}$$

By Theorem 2.2 of Bruss and Robertson (1991), it holds that, for any $\alpha \in (0,1)$,

$$\lim_{n \to \infty} \frac{1}{n} E[M_n^*(n \cdot c(\alpha))] = \alpha$$

where $c(\alpha) = \int_0^{F^{-1}(\alpha)} x dF(x)$. By setting $\alpha = \beta$, we have

$$\lim_{n \to \infty} \frac{1}{n} E[M_n^*(n \cdot c(\beta))] = \beta$$

Comparing the above equality with (71), we concludes that the capacity level $n \cdot c^*$ is asymptotically optimal as $n \to \infty$, where $c^* = c(\beta)$. Note that the following optimization is a concave maximization problem,

$$\max_{\xi}\{\xi - \frac{1}{\beta} \mathrm{E}_{\tilde{D}}[(\xi - \tilde{D})^+]\} \tag{72}$$

and the maximum is achieved by setting the derivative of the objective function with respect to $\xi$ to be 0. Denote by

$$Z(\xi) := \xi - \frac{1}{\beta} \mathrm{E}_{\tilde{D}}[(\xi - \tilde{D})^+]$$

the function within the min operation in (72), and denote by $\xi^* = \mathrm{argmax} Z(\xi)$. We have

$$\frac{\partial}{\partial \xi} Z(\xi^*) = 1 - \frac{1}{\beta} F(\xi^*) = 0,$$

which implies that $\xi^* = F^{-1}(\beta)$. Thus, we have

$$\max_{\xi}\{\xi - \frac{1}{\beta} \mathrm{E}_{\tilde{D}}[(\xi - \tilde{D})^+]\} = \frac{1}{\beta} \cdot \int_0^{F^{-1}(\beta)} x dF(x) = \frac{1}{\beta} \cdot c^*$$

We then bound $c^*$ and prove (68). Note that

$$\beta\mu - c^* = \beta \cdot \int_0^\infty x dF(x) - \int_0^{\xi^*} x dF(x) = \beta \cdot \int_{\xi^*}^\infty x dF(x) - (1-\beta) \cdot \int_0^{\xi^*} x dF(x)$$

$$\ge \beta \cdot \int_{\xi^*}^\infty \xi^* dF(x) - (1-\beta) \cdot \int_0^{\xi^*} \xi^* dF(x) = \beta(1-\beta)\xi^* - (1-\beta)\beta\xi^* = 0$$

Then, we have $c^* \le \beta\mu$. Also from Cauchy inequality, we have that

$$(1-\beta) \cdot \int_{\xi^*}^\infty x^2 dF(x) \ge (\int_{\xi^*}^\infty x dF(x))^2$$

which implies

$$(1-\beta)\cdot[\mu^2+\sigma^2-\int_0^{\xi^*}x^2dF(x)] \ge (\mu-\int_0^{\xi^*}xdF(x))^2$$

Thus, we have

$$(1-\beta)\cdot(\mu^2+\sigma^2) \ge (\mu-\int_0^{\xi^*}xdF(x))^2+(1-\beta)\cdot\int_0^{\xi^*}x^2dF(x) \ge (\mu-\int_0^{\xi^*}xdF(x))^2+\frac{1-\beta}{\beta}\cdot(\int_0^{\xi^*}xdF(x))^2$$

By arranging terms in the above inequality, we have

$$(\beta\mu-\int_0^{\xi^*}xdF(x))^2 \le \sigma^2\beta(1-\beta)$$

Thus, it holds that

$$c^* = \int_0^{\xi^*}xdF(x) \ge \beta\mu-\sigma\cdot\sqrt{\beta(1-\beta)}$$

which completes our proof. $\square$

## Appendix J: Proof of Proposition 1

Suppose the capacity level $c$, together with a rationing policy $\tilde{\phi}^{(1)}$, is feasible for differentiated service levels $\beta = (\beta_1, \beta_2, \ldots, \beta_n)$. We say a policy $\tilde{\phi}'$ is the one-step rotation of policy $\tilde{\phi}$ if $s_j(\tilde{\phi}', c, \mathbf{D}) = s_{j+1}(\tilde{\phi}, c, \mathbf{D}), j = 1, 2, \ldots, n-1$, and $s_n(\tilde{\phi}', c, \mathbf{D}) = s_1(\tilde{\phi}, c, \mathbf{D})$. Let $\tilde{\phi}^{(k+1)}$ be the one-step rotation of $\tilde{\phi}^{(k)}$ for all $k$ from 1 to $n-1$. Since demand distributions are i.i.d., $\tilde{\phi}^{(k)}$ must be feasible for service levels $(\beta_k, \beta_{k+1}, \ldots, \beta_n, \beta_1, \beta_2, \ldots, \beta_{k-1})$. Thus the randomized policy that chooses $\tilde{\phi}^{(k)}$ with probability $\frac{1}{n}$, for each $k = 1, 2, \ldots, n$, can achieve a service level $\hat{\beta} = \frac{1}{n}\sum_{j=1}^n \beta_j$ for each customer $j \in \mathcal{N}$. Henceforth, the optimal capacity level with uniform target service level $\hat{\beta} = \frac{1}{n}\sum_{j\in\mathcal{N}}\beta_j$ for all customer $j$ is less than or equal to $c$. $\square$