

Optimal and Near-optimal Policies of Sequential Decision Problems under Uncertainty

by

Xinchang Xie

Business Administration
Duke University

Date: _____
Approved: _____

Alessandro Arlotto, Advisor

Alexandre Belloni, Co-Advisor

Yehua Wei

Shankar Bhamidi

Dissertation submitted in partial fulfillment of the
requirements for the degree of Doctor of Philosophy
in the Department of Business Administration
in the Graduate School of
Duke University

2020

ABSTRACT

Optimal and Near-optimal Policies of Sequential Decision Problems
under Uncertainty

by

Xinchang Xie

Business Administration
Duke University

Date: _____
Approved: _____

Alessandro Arlotto, Advisor

Alexandre Belloni, Co-Advisor

Yehua Wei

Shankar Bhamidi

An abstract of a dissertation submitted in partial fulfillment of the
requirements for the degree of Doctor of Philosophy
in the Department of Business Administration
in the Graduate School of
Duke University

2020

Abstract

In this dissertation, we study stochastic sequential decision problems with a discrete finite time horizon. We conduct substantial analyses on both optimal and near-optimal policies of such problems, and focus on the *total reward* collected by different sequential policies from different perspectives. The total reward collected by any sequential policy is *a priori* random, and for many of the policies we study, we are interested in: (i) the expected total reward they collect, (ii) the variance of their total reward, (iii) the limiting distribution of their total reward, and (iv) confidence intervals for their total reward. With these goals in mind, we consider different sequential decision problems that include the dynamic and stochastic knapsack problem with equal rewards, the problem of sequential monotone subsequence selection from a random sample, and a dynamic inventory control problem. For the dynamic and stochastic knapsack problem with equal rewards, we propose an adaptive heuristic policy and prove that its regret—the expected gap between the total reward collected by the proposed heuristic policy and that collected by the best offline solution—is at most logarithmically in the problem size. Furthermore, we characterize the variance asymptotics as well as the limiting distribution of the total reward collected by such heuristic. We show that the total reward collected by our heuristic policy shares the same variance asymptotics and the same limiting distribution with that collected by the optimal sequential policy. This is in contrast with the performance of other asymptotically optimal heuristics whose total rewards have larger regrets, larger variances, and different limiting distributions. We also discuss the equivalence between the problem of sequential monotone subsequence selection from a random sample and a special instance of the dynamic and stochastic knapsack problem with equal rewards. Such equivalence enables us to apply the above mentioned analyses and results to the sequential monotone subsequence selection problem, and establish a number of results of independent interest. Lastly, we study multiple finite horizon sequential decision problems that are faced in parallel by self-interested decision makers. The uncertainties they face could be correlated in an arbitrary unknown fashion. Our goal is to construct simultaneous confidence intervals for the total rewards collected by each decision maker, through the observations of the rewards they collect along the time horizon. We propose a data-driven bootstrap procedure that is asymptotically valid. The key feature of such procedure is that the number of decision makers is allowed to be much larger than the length of the time horizon. This

corresponds to the high-dimensional regime in the statistics literature in which the dimension of the estimator is much larger than sample size. As a byproduct, such bootstrap procedure can be used to simultaneously test whether each of the decision makers is implementing certain specified policy. We use a dynamic inventory control problem and the sequential monotone subsequence selection problem to illustrate the effectiveness of our framework.

Acknowledgments

The Ph.D. process has been a challenging yet enjoyable journey for me, and I could not have completed it without my co-advisors Prof. Alessandro Arlotto and Prof. Alex Belloni.

It is my great fortune to have Alessandro as my co-advisor, and I could not have made it to accomplish this milestone without his guidance and support. Through our weekly meetings since my second year of the program, I have learned from him not only how to conduct high standard research, but also how to (and how important it is to) articulate research results in papers and presentations. He has also been guiding me with patience and encouragement, through which I grew—step by step—as a researcher and as a person.

I am also deeply thankful to Alex for serving my other co-advisor. With his broad research interests, knowledge, and expertise, Alex has taught me how to connect different fields to find out interesting research questions. He has been truly supportive and would always make time out of his busy schedule to address my questions. He has also been leading by example and teaching me the importance of hardworking to succeed.

My sincere thanks also go to the other two professors on my committee: Prof. Yehua Wei, who was my mentor during my first year and has been supporting me since then, and Prof. Shankar Bhamidi, who kindly served on my committee as an external member and has provided valuable comments and perspectives for my candidacy exam, proposal exam, and my final defense. In addition, I want to thank several other faculty members at Fuqua: Peng Sun, David Brown, Bob Nau, Bob Winkler, and Bora Keskin. Through their courses, I deepened my knowledge and learned tools that have been helpful for my research.

Besides all the help I got from the faculty, family and friends have been supporting me since the very beginning. My deepest appreciation is for my parents, who have been supporting me unconditionally all the time. I would also like to thank my dear friends within the Decision Sciences area and the Operations Management area at Fuqua: Levi DeValve, Huseyin Gurkan, Mingliu Chen, Chen Chen, Yunke Mai, and Soudipta Chakraborty—with whom I took classes with, and discussed research problems through the years. More importantly, I am grateful for the many Fuqua Fridays that we spent together. I am also thankful for my other friends across different departments at

Duke: Yanyou Chen, Joy Tong, Anyi Ma, Wenxi Liao, Ye Jin, Xuyan Ru, David Hall, and Honggi Lee. It is you guys who made this journey more enjoyable and unforgettable.

Finally, I am gratefully acknowledging partial financial support from the National Science Foundation (Grant No. 1553274).

Contents

Abstract	iv
Acknowledgments	vi
List of Figures	xi
List of Tables	xii
1 Introduction	1
2 Logarithmic Regret in the Dynamic and Stochastic Knapsack Problem with Equal Rewards	4
2.1 Literature Review: Knapsack Problems and Approximations	10
2.2 A Prophet Upper Bound	12
2.3 The Reoptimized Policy and Its Value Function	17
2.4 On the Typical Class	19
2.5 A Logarithmic Regret Bound	22
2.5.1 Preliminary Observations	23
2.5.2 Analysis of Residuals	27
2.6 Numerical Experiments	31
2.7 On Weight Distributions with Multiple Types	34
2.8 Concluding Remarks	37
3 An Adaptive $O(\log n)$-Optimal Policy for the Sequential Selection of a Monotone Subsequence from a Random Sample	38
3.1 Policy $\hat{\pi}$ and Its Value Function	41
3.2 A Refined Prophet Upper Bound	43
3.3 Equivalence with the Dynamic and Stochastic Knapsack Problem with Equal Rewards	45
3.4 Connections and Observations	50
4 Sequential Policies and the Distribution of Their Total Rewards in Dynamic and Stochastic Knapsack Problems	52

4.1	Dynamic and Stochastic Knapsack Problem with Unitary Rewards: A Distributional Perspective	55
4.1.1	The Adaptive Heuristic Revisited	56
4.1.2	Main Result	60
4.1.3	Properties of the Regular Class	62
4.2	Non-asymptotic Derivative Bounds	69
4.3	Variance Asymptotics with Remainder Bounds	85
4.3.1	Doob's Martingale and Conditional Variance of Martingale Differences	85
4.3.2	Conditional Variance Bounds	87
4.4	Martingale Central Limit Theorem	96
4.5	Concluding Remarks	100
5	Data-driven Monitoring the Implementation of Policies across Many Markov Decision Problems	101
5.1	Literature Review	104
5.2	A General MDP Framework	106
5.2.1	Hypothesis Testing and the Construction of Simultaneous Confidence Intervals	107
5.2.2	Martingale Representation	107
5.3	MDPs with Regenerative Property	108
5.3.1	Test Statistic and the Block Multiplier Bootstrap	109
5.3.2	Asymptotic Validity	111
5.4	MDPs with Absorbing States	111
5.4.1	Test Statistics and the Multiplier Bootstrap	112
5.4.2	Asymptotic Validity	113
5.5	Application: Testing Stationary Inventory Control Policies	114
5.5.1	Problem Setup	114
5.5.2	Regeneration Time with Exponential Tail	117
5.5.3	Analysis on the Martingale Differences	122

5.6	Application: Online Monotone Subsequence Selection from a Random Sample	125
5.6.1	Problem Description and Main Result	126
5.6.2	Control the Covariance Matrix	128
5.6.3	Verification of Other Conditions	142
5.7	Concluding Remarks	145
6	Conclusion	146
	Bibliography	147

List of Figures

2.1	Gap between the prophet upper bound and offline sort for three weight distributions.	31
2.2	Value functions and scaled regret bounds for three weight distributions .	32

List of Tables

4.1	Asymptotic performance comparison as $n \rightarrow \infty$ among three policies.	55
-----	---	----

Chapter 1

Introduction

In this dissertation, we study finite-horizon sequential decision problems under uncertainty, and analyze their optimal and near-optimal policies. Sequential decision problems are widely used in operations management and operations research; their versatility can be easily seen through a wide range of applications. The finite horizon formulation of such problems can be typically described as follows. In each time period, a decision maker (referred to as *she*) sees the realization of an exogenous stochastic process and has to decide what action to take based on the available information and on her past actions. The realized uncertainty and the chosen action then determine a reward she collects during that time period. The objective is typically to maximize the expected total collected reward. The optimal policies of such problems are usually difficult to analyze as that they are both time- and state-dependent. While one can in principle analyze such dependencies by solving the problems to optimality through the associated Bellman equations, such solutions are often lack of closed-form expressions, and their computations suffer the infamous *curse of dimensionality* (Bellman, 1957; Powell, 2011).

As a substitute, decision makers often seek for near-optimal heuristics with provable performance guarantees. Such heuristics are usually evaluated in terms of *optimality gap* or *regret*. Optimality gap is the expected performance gap between a heuristic and the optimal sequential policy. Instead, regret is the expected performance gap between a heuristic and the optimal offline policy. Closely related to these two concepts is the notion of *asymptotic optimality*. A heuristic sequential policy is said to be asymptotically optimal if the ratio between its optimality gap and the expected performance of the optimal sequential policy goes to zero as the problem size goes to infinity. In this dissertation, we will evaluate different policies for certain problems in terms of these metrics.

In addition to these metrics that are based on the expected performance, we also study sequential decision problems from a distributional standpoint. We analyze the variance of the total reward, which is critical to differentiate policies that have asymptotically equivalent means. We also analyze of the limiting distribution of the total reward, which is helpful for constructing confidence intervals for the total reward. All such analyses have both theoretical and practical benefits.

From a theoretical standpoint, such analyses would help us further differentiate different policies within the class of asymptotically optimal policies. From a practical standpoint, the construction of confidence intervals could be applied to compute the value at risk of different policies and to test whether decision makers are implementing certain policies through the observations of the rewards they collect over time.

The sequential decision problems we will cover in this dissertation include the dynamic and stochastic knapsack problem with equal rewards, the problem of sequential monotone subsequence selection from a random sample, and a dynamic inventory control problem. In Chapter 2, we study the dynamic and stochastic knapsack problem with equal rewards. Such a formulation of knapsack problem was studied by Coffman et al. (1987) and Kleywegt and Papastavrou (1998, Section 5). For this problem, we propose and study a simple adaptive heuristic with closed-form expression, and prove that its regret—the expected performance difference between such heuristic and the offline optimal solution—grows at most logarithmically in the problem size. In Chapter 3, we discuss how the logarithmic regret bound of Chapter 2 carries over to the equivalent problem of sequential selection of an increasing subsequence from a random sample.

Then in Chapter 4, we go back to the dynamic and stochastic knapsack problem with equal rewards and offer a distributional analysis. Specifically, we focus on the higher-order performance and limiting behavior of the heuristic proposed in Chapter 2, and prove that it is not only asymptotically optimal, but also that it collects a total reward that has the same variance asymptotics and the same limiting distribution as that collected by the optimal sequential policy. We note that the total reward collected by other asymptotically optimal heuristic studied in the literature may have different (larger) variance and a different limiting distribution than that collected by the optimal sequential policy.

When one has the available limiting distribution of the total reward collected by a policy for a sequential decision problem, then it is easy to construct confidence intervals for the total reward, or to conduct hypothesis testing of whether the decision maker is implementing such a policy. In Chapter 5, we consider such two statistical tasks when there are many sequential decision problems faced by self-interested decision makers and the uncertainties they face are correlated in an arbitrary unknown fashion. To this end, we propose a bootstrap based procedure for constructing simultaneous confidence intervals for the total rewards collected by each decision maker. As a key

feature, such procedure allows for the number of decision makers being comparable to or even much larger than the length of the time horizon. This regime is often referred to as *high dimensional*; the the number of decision makers can be viewed as the data's dimension and the length of the time horizon can be interpreted as sample size. Lastly, we use a dynamic inventory control problem and the sequential monotone subsequence selection problem to illustrate the effectiveness of the framework.

Chapter 2

Logarithmic Regret in the Dynamic and Stochastic Knapsack Problem with Equal Rewards

The knapsack problem is one of the classic problems in operations research. It arises in resource allocation, and it counts numerous applications in auctions, logistics, portfolio optimization, scheduling, and transportation among others (cf. Martello and Toth, 1990; Kellerer et al., 2004). In its dynamic and stochastic formulation (see, e.g. Papastavrou et al., 1996; Kleywegt and Papastavrou, 1998, 2001) a decision maker (referred to as *she*) is given a knapsack with finite capacity $0 \leq c < \infty$ and is sequentially presented with items arriving over a time horizon with n discrete time periods, indexed by $i \in [n] \equiv \{1, 2, \dots, n\}$. In each period $i \in [n]$, an item arrives with probability p , its weight-reward pair $(\mathfrak{W}_i, \mathfrak{R}_i)$ is revealed, and the decision maker needs to decide whether to include the arriving item in the knapsack or to reject it forever. Here, the weight \mathfrak{W}_i represents the amount of knapsack capacity that the item arriving in period i consumes if the decision maker chooses to include it in the knapsack, and the reward \mathfrak{R}_i represents what the decision maker collects upon inclusion. The pairs $(\mathfrak{W}_i, \mathfrak{R}_i)$, $i \in [n]$, are independent and with common, known, bivariate distribution supported on the nonnegative orthant.

By imposing different assumptions on the weight-reward distribution, one recovers knapsack instances of independent interest. For instance, in the problem of real-time uniprocessor scheduling under conditions of overload (see, e.g., Baruah et al., 1994) a decision maker wants to maximize the number of jobs that are processed on a single machine by a fixed deadline. In this context, the deadline is the knapsack capacity and jobs correspond to items. Their rewards are all equal to one, and their durations correspond to the item weights. This scheduling application motivates the model in this chapter. We assume that the rewards are deterministic and all equal¹ to $r > 0$, and the weights are independent random variables with common continuous distribution F . We model item arrivals by considering a Bernoulli process $\mathfrak{B}_1, \mathfrak{B}_2, \dots, \mathfrak{B}_n$ that is independent of everything else, and that is given by a sequence of independent Bernoulli random variables with success probability

This chapter is written under the supervision of Prof. Alessandro Arlotto. The results presented here are also in the joint paper Arlotto and Xie (2020a), forthcoming at *Stochastic Systems*.

¹This also covers random rewards with common distribution that are revealed only after the inclusion decision.

p . We then equivalently redefine the weight distribution so that a no arrival corresponds to the arrival of an item with arbitrarily large weight. That is, we assume that an item arrives in each period $i \in [n]$ and that its weight is given by the random variable W_i defined by

$$W_i = \begin{cases} +\infty & \text{if } \mathfrak{B}_i = 0 \\ \mathfrak{W}_i & \text{if } \mathfrak{B}_i = 1. \end{cases}$$

We say that a policy π is *feasible* if the sum of the weights of the items selected by π does not exceed the knapsack capacity c , and we say that the policy is *online* (or *sequential*) if the decision to select item i with weight W_i depends only on the information available up to and including time i . We then let $\Pi(n, c, p)$ be the set of feasible online policies, and we compare the performance of the best online policy to that of a prophet who has full (or *offline*) knowledge of the weights W_1, W_2, \dots, W_n before making any selection. Under some mild technical conditions on the weight distribution F , we prove that the regret—the expected gap between the performance of the best online policy and its offline counterpart—is bounded by the logarithm of n . Our proof is constructive. We propose a reoptimized heuristic that exhibits logarithmic regret. The heuristic is based on resolving some related optimization problem at any given time $i \in [n]$ by using the current—rather than the initial—level of remaining capacity as constraint. The solution of this optimization problem provides us with a state- and time-dependent threshold that mimics that of the optimal online policy.

If all of the weights W_1, W_2, \dots, W_n are revealed to the decision maker before she makes any selection, then her choice is obvious. To maximize the total reward she collects, she just sorts the items according to their weights and selects them starting from the smallest weight and continuing until the knapsack capacity is exhausted. Formally, if $W_{(1,n)} \leq W_{(2,n)} \leq \dots \leq W_{(n,n)}$ are the order statistics of W_1, W_2, \dots, W_n , then the maximal reward $R_n^*(c, p, r)$ that the decision maker collects is given by

$$R_n^*(c, p, r) = \max \left\{ rm : m \in \{0, 1, \dots, n\} \text{ and } \sum_{\ell=1}^m W_{(\ell,n)} \leq c \right\}. \quad (2.1)$$

Here we compare the total reward of the offline-sort algorithm (2.1), $R_n^*(c, p, r)$, with that of an online feasible policy $\hat{\pi} \in \Pi(n, c, p)$ that is based on a sequence of reoptimized time- and state-dependent threshold functions $\hat{h}_n, \hat{h}_{n-1}, \dots, \hat{h}_1$. If the current level of remaining capacity is x and the weight of item i is about to be revealed, then the decision maker computes the threshold $\hat{h}_{n-i+1} : [0, \infty) \rightarrow [0, \infty)$ such that $\hat{h}_{n-i+1}(x) \leq x$, and she selects item i if and only if the weight

$W_i \leq \hat{h}_{n-i+1}(x)$. Thus if $\hat{X}_0 = c$ and for $i \in [n]$ one defines the remaining capacity process \hat{X}_i recursively by

$$\hat{X}_i = \begin{cases} \hat{X}_{i-1} & \text{if } W_i > \hat{h}_{n-i+1}(\hat{X}_{i-1}) \\ \hat{X}_{i-1} - W_i & \text{if } W_i \leq \hat{h}_{n-i+1}(\hat{X}_{i-1}), \end{cases}$$

then the total reward collected by the reoptimized policy $\hat{\pi}$ can be written as

$$R_n^{\hat{\pi}}(c, p, r) = \sum_{i=1}^n r \mathbb{1} \left\{ W_i \leq \hat{h}_{n-i+1}(\hat{X}_{i-1}) \right\}.$$

The random variables $R_n^*(c, p, r)$ and $R_n^{\hat{\pi}}(c, p, r)$ crucially depend on the weight distribution F . This dependence is mostly expressed through a *consumption function* $\epsilon_{kp} : [0, \infty) \rightarrow [0, \infty]$ that is defined for $p \in (0, 1]$ and for all $1 \leq k < \infty$ by

$$\epsilon_{kp}(x) = \sup \left\{ \epsilon \in [0, \infty) : \int_0^\epsilon w dF(w) \leq \frac{x}{kp} \right\}. \quad (2.2)$$

The consumption function depends on two quantities. The argument x that denotes the current level of remaining capacity of the knapsack, and the index kp that refers to the expected number of items with F -distributed weights (or *arrivals*) that are yet to be presented to the decision maker. Furthermore, the function $\epsilon_{kp}(x)$ is always well defined. If $\mu = \mathbb{E}[\mathfrak{W}_1] = \mathbb{E}[W_1 \mid B_1 = 1]$ and $kp\mu < x < \infty$ then $\epsilon_{kp}(x) = +\infty$. Otherwise, the value $\epsilon_{kp}(x)$ satisfies the integral representation

$$\int_0^{\epsilon_{kp}(x)} w dF(w) = \frac{x}{kp} \quad \text{for all } x \in [0, kp\mu]. \quad (2.3)$$

The representation (2.3) offers an important insight regarding the role of the consumption function $\epsilon_{kp}(x)$. The integral on the left-hand side is the *expected reduction* in the remaining capacity of the knapsack when the current level of remaining capacity is equal to x , and the decision maker selects an item with weight smaller than $\epsilon_{kp}(x)$. The function $\epsilon_{kp}(x)$ is then defined so that the expected reduction in capacity is equal to the ratio of the current capacity, x , to the expected number of remaining arrivals, kp . That is, the threshold $\epsilon_{kp}(x)$ is constructed so that—in expectation—the available capacity is spread equally over the remaining arrivals.

As we will see shortly, the threshold $\epsilon_{kp}(x)$ drives most of the estimates in this chapter and, together with the continuity of the weight distribution F , it immediately provides us with an easy upper bound for $\mathbb{E}[R_n^*(c, p, r)]$. The same threshold together with some mild regularity conditions

on the weight distribution F also drives the lower bound for $\mathbb{E}[R_n^{\hat{\pi}}(c, p, r)]$. The class of weight distributions we consider for the lower bound is characterized in the next definition.

Definition 2.1 (Typical class of distributions with continuous density). We say that a non-negative distribution F with continuous density function f belongs to the *typical class* if for some $\bar{w} > 0$, the following two conditions hold.

(i) BEHAVIOR AT ZERO. There are $0 < \lambda < 1$ and $0 < \gamma < 1$ such that

$$\frac{F(\lambda w)}{F(w)} \leq \gamma < 1 \quad \text{for all } w \in (0, \bar{w}). \quad (2.4)$$

(ii) MONOTONICITY. The map $w \mapsto w^3 f(w)$ is non-decreasing on $(0, \bar{w})$. That is,

$$w_1^3 f(w_1) \leq w_2^3 f(w_2) \quad \text{for all } 0 < w_1 \leq w_2 < \bar{w}. \quad (2.5)$$

The class of typical distributions is wide enough to include most well-known non-negative distributions. In Section 2.4, we provide specific examples as well as class properties, but for now we emphasize that the breadth of the typical class comes from the role of the distribution-dependent parameter $\bar{w} > 0$. Conditions (2.4) and (2.5) need only to hold near zero—or, more precisely, on $(0, \bar{w})$ —and not on the full support of the weight distribution or on the whole capacity interval $[0, c]$. In fact, for many distributions the parameter \bar{w} for which (2.4) and (2.5) hold is much smaller than the minimum between the initial capacity and the supremum of the support.

The main results of this chapter are gathered in the theorem below. First, we provide an upper bound for $\mathbb{E}[R_n^*(c, p, r)]$ that holds for any continuous distribution F . Then, we turn to distributions that belong to the typical class, and we prove that there is a matching lower bound. As a by-product of our analysis, we establish that the regret is, at most, $O(\log n)$ as $n \rightarrow \infty$.² While our theoretical result provides only a regret bound, related results and the numerical experiments of Section 2.6 tell us that the regret bound is actually of the correct order.

Theorem 2.1 (Logarithmic regret bound). *Consider a knapsack problem with capacity $0 \leq c < \infty$ and with items that arrive over $1 \leq n < \infty$ periods according to a Bernoulli process with arrival probability $p \in (0, 1]$. If the items have rewards equal to r and weights with continuous distribution F , then*

$$\max_{\pi \in \Pi(n, c, p)} \mathbb{E}[R_n^\pi(c, p, r)] \leq \mathbb{E}[R_n^*(c, p, r)] \leq nprF(\epsilon_{np}(c)).$$

²Throughout this chapter, the function \log denotes the natural logarithm.

Furthermore, there is a feasible online policy $\hat{\pi} \in \Pi(n, c, p)$ such that if the weights are independent and their distribution F belongs to the typical class then there is a constant $1 < M < \infty$ depending only on F , p , and r for which

$$nprF(\epsilon_{np}(c)) - M(1 + \log n) \leq \mathbb{E} \left[R_n^{\hat{\pi}}(c, p, r) \right] \leq \max_{\pi \in \Pi(n, c, p)} \mathbb{E} [R_n^{\pi}(c, p, r)].$$

In turn, if the weights are independent and the distribution F belongs to the typical class, then we have the regret bound

$$\mathbb{E} [R_n^*(c, p, r)] - \max_{\pi \in \Pi(n, c, p)} \mathbb{E} [R_n^{\pi}(c, p, r)] \leq \mathbb{E} [R_n^*(c, p, r)] - \mathbb{E} [R_n^{\hat{\pi}}(c, p, r)] \leq M(1 + \log n).$$

The special case with deterministic arrivals and unitary rewards has been extensively studied in the literature. The upper bound $\mathbb{E} [R_n^*(c, 1, 1)] \leq nF(\epsilon_n(c))$ was first proved by Bruss and Robertson (1991). Here, we provide a generalization that is based on a relaxation of some appropriate optimization problem. The solution to this relaxation is the basis for constructing the reoptimized heuristics $\hat{\pi}$. The lower bound $\mathbb{E} [R_n^{\hat{\pi}}(c, p, r)] \geq nprF(\epsilon_{np}(c)) - O(\log n)$ as $n \rightarrow \infty$ is essentially new, and it substantially improves on existing estimates. The best results to date for general weight distribution F are due to Rhee and Talagrand (1991) who study a non-adaptive heuristic and prove that

$$nF(\epsilon_n(c)) \left\{ 1 - \left[\frac{\epsilon_n(c)}{c} \right]^{1/2} - \frac{\epsilon_n(c)}{c} \right\} \leq \max_{\pi \in \Pi(n, c, 1)} \mathbb{E} [R_n^{\pi}(c, 1, 1)] \quad \text{for all } n \geq 1. \quad (2.6)$$

For instance, if $F(x) = \sqrt{x}$ for $x \in (0, 1)$ then the lower bound (2.6) implies an upper bound for the regret that is $O(n^{1/3})$ as $n \rightarrow \infty$. Similarly, if $F(x) = x^2$ for $x \in (0, 1)$ then the same lower bound gives us a regret upper bound that behaves like $O(n^{1/6})$ as $n \rightarrow \infty$.

A case that deserves special attention is when F is the uniform distribution on the unit interval, the reward $r = 1$, and the initial capacity $c = 1$. In this context, the Rhee and Talagrand (1991) lower bound provides us with a regret upper bound that behaves like $O(n^{1/4})$ as $n \rightarrow \infty$, but better bounds are available in the literature. This special dynamic and stochastic knapsack problem is in fact equivalent to the problem of the sequential selection of a monotone decreasing subsequence from a sample of n independent observation with the uniform distribution on the unit interval (cf. Samuels and Steele, 1981). The equivalence was first observed by Coffman et al. (1987, pp. 457–458), and it can be established by observing that the Bellman equations for the two problems are the

same after a change of variable. Informally, if the number of remaining periods is the same in both problems and the current capacity of the knapsack is equal to the last selected subsequence element, then the largest weight that is optimal for inclusion is equal to the maximum amount the decision maker is willing to go down in optimally selecting a new subsequence element. Since the weights as well as the subsequence elements are both uniformly distributed on the unit interval, these two actions happen with the same probability. For this subsequence-selection problem, Arlotto et al. (2015, 2018) prove that the expected performance ν_n^* of the best online policy satisfies the estimate $\nu_n^* = \sqrt{2n} - O(\log n)$ as $n \rightarrow \infty$. The equivalence between the two problems, however, holds *only* for uniform weights. As Theorem 2.1 suggests, the weight distribution F plays a crucial role in the estimates for the dynamic and stochastic knapsack problem with equal rewards. Instead, the monotone subsequence problem is distribution invariant, and one can consider uniformly distributed subsequence elements without loss of generality. More importantly, Seksenbayev (2018) and Gnedin and Seksenbayev (2019) characterize the second order asymptotic expansion of ν_n^* and establish that $\nu_n^* = \sqrt{2n} - \frac{1}{12} \log n + O(1)$ as $n \rightarrow \infty$. This remarkable result tells us that our regret bound is order tight, and that no online algorithm can—at this level of generality—be within $O(1)$ of offline sort.

Organization of the Chapter

The rest of this chapter is organized as follows. In Section 2.1, we review the related literature. In Section 2.2, we prove the prophet upper bound $\mathbb{E}[R_n^*(c, p, r)] \leq nprF(\epsilon_{np}(c))$ by showing that the offline-sort algorithm (2.1) can be reinterpreted as a parsimonious threshold policy and by solving a relaxation of some related optimization problem. This solution then guides us in the construction of policy $\hat{\pi}$ that is presented in Section 2.3. In Section 2.4, we discuss the generality of the typical class of distributions, and we derive some properties that we then use—in Section 2.5—to prove that the reoptimized policy $\hat{\pi}$ exhibits logarithmic regret. In Section 2.6, we present numerical experiments that provide further insights into our regret bound, while in Section 2.7 we discuss weight distributions with multiple types. Finally, in Section 2.8 we make closing remarks and underscore some open problems.

2.1 Literature Review: Knapsack Problems and Approximations

Knapsack problems uniquely combine simple formulations, non-trivial mathematical analyses, and relevance in several application-driven domains. As such, different knapsack problems have been considered in the literature, and a lot of effort has been devoted to the development of (near-) optimal policies. Most of the differences that have been accounted for concern the item arrival process (*static* versus *dynamic*), the probabilistic assumptions on the weight-reward pairs (deterministic and/or stochastic), and the objective of the decision maker (reward maximization, target achievement, etc.).

For instance, in the early formulation of Dantzig (1957), we have a *static* model with a finite number of items that are all available before any decision is made and have deterministic weights and deterministic rewards. The decision maker then seeks to find a maximum-reward subset of these items with total weight that does not exceed a capacity constraint. Following this classic formulation, researchers have considered several *static* knapsack instances with randomness in the weights and/or in the rewards. While studying a scheduling problem, Derman et al. (1978) studied a *static* and stochastic knapsack problem with items that belong to different categories. Items that belong to the same category have common deterministic rewards and independent, exponentially distributed weights with category-dependent parameter. The decision maker then seeks to maximize total expected rewards when the realized weights are revealed only after each item is included in the knapsack. The authors prove that the greedy policy based on reward-to-mean-weight ratios is optimal. Analogous *static* and stochastic knapsack problems have been considered by several authors, including Dean et al. (2004, 2005, 2008), Bhargat et al. (2011), Li and Yuan (2013), Blado et al. (2016), Ma (2018), Blado and Toriello (2019), and Balseiro and Brown (2019). Gupta et al. (2011) and Merzifonluoglu et al. (2012) follow along similar lines, but consider both random weights and random rewards. Most notably, Dean et al. (2004, 2005, 2008) study a *static* and stochastic knapsack problem with deterministic rewards and independent random weights with arbitrary distributions that are realized only upon insertion in the knapsack. They construct a polynomial time adaptive policy that is within a constant multiplicative gap, and they compare the performance of adaptive and non-adaptive policies. Their work is particularly relevant to us as it is among the first ones to assess the benefits of adaptivity.

Static stochastic knapsack problems have also been studied under different optimization objectives. For instance, there is a stream of related literature that considers *static* stochastic knapsack problems (typically with deterministic weights and random rewards) in which the objective is to maximize the probability that the total reward will achieve a certain given target. (See, e.g., Henig, 1990; Carraway et al., 1993; Ilhan et al., 2011, among others.)

Alongside the static knapsack problems mentioned thus far there are several *dynamic* models in which items arrive over time and their weight-reward pairs are revealed to the decision maker who irrevocably decides on inclusion in the knapsack as soon as each item arrives and without seeing the weights and/or the rewards of future items. *Dynamic* and stochastic knapsack problems are widespread. For instance, if one assumes that the weights are all equal to one and that the rewards are random, then one recovers the multi-secretary problem (see, e.g. Cayley, 1875; Moser, 1956; Kleinberg, 2005). For this problem, Arlotto and Gurvich (2019) prove that if the reward distribution is discrete, then the regret is uniformly bounded in the number of items and the knapsack capacity. Similarly, if one assumes that the rewards are all equal to one and that the weights are random, then one finds an instance of the single-machine scheduling problem of Baruah et al. (1994) that motivates this chapter. Finally, when both the weights and the rewards are random, one recovers—among others—the sequential investment problems of Derman et al. (1975) and Prastacos (1983), or the multi-secretary problem of Nakai (1986) which allows for an unknown number of applicants in each period. When both the weights and the rewards are random, few regret bounds are available. A notable exception is the work of Marchetti-Spaccamela and Vercellis (1995) who prove a $O(\log^{3/2} n)$ regret bound when both the weights and the rewards are independent and uniformly distributed on the unit interval, and the knapsack capacity is proportional to the number of periods. For the same formulation, Lueker (1998) improves Marchetti-Spaccamela and Vercellis’s result to $O(\log n)$ and shows that it is best possible.

Multi-dimensional generalizations of the *dynamic* and stochastic knapsack problem have found several applications in revenue management and resource allocation. In the network revenue management problem, heterogeneous customers belonging to different classes arrive sequentially over time, request a product, and offer a price. If the request is accepted, then a collection of resources that constitute the product is depleted, and the offered price is earned. Otherwise the resource capacities remain unchanged and the offered price is lost (cf. Gallego and van Ryzin, 1997; Talluri and van Ryzin, 2004). The solution of the network revenue management problem is famously diffi-

cult, and scholars have studied several non-adaptive as well as adaptive heuristics and proved regret bounds. A classic non-adaptive approximation scheme based on a deterministic linear-programming relaxation was studied by Gallego and van Ryzin (1994, 1997). In contrast, adaptive policies have been considered by allowing for periodic reoptimization. Despite a few specific negative results by Cooper (2002), Chen and Homem-de Mello (2010), and Jasin and Kumar (2013), there are ways to construct reoptimized policies that perform well. For instance, Reiman and Wang (2008) propose a probabilistic allocation rule that works well with one reoptimization instance. Jasin and Kumar (2012) and Wu et al. (2015) consider a probabilistic allocation rule that is based on reoptimizing in every period and show that it exhibits uniformly bounded regret provided that the optimal solution to the original deterministic linear programming relaxation is non-degenerate. Bumpensanti and Wang (2020) and Vera and Banerjee (2019) prove that the uniform regret bound holds in general, without the non-degeneracy assumption.

2.2 A Prophet Upper Bound

The performance of any online algorithm is bounded above by the full-information (or offline) sort. If the decision maker knows all of the weights W_1, W_2, \dots, W_n before making any decision, then the total reward she collects is the largest number rm such that the sum of the smallest m realizations does not exceed the capacity constraint. That is, if $W_{(1,n)} \leq W_{(2,n)} \leq \dots \leq W_{(n,n)}$ are the order statistics of $\mathcal{W} \equiv \{W_1, W_2, \dots, W_n\}$, then the total reward $R_n^*(c, p, r)$ of offline selections when the initial knapsack capacity is c and the arrival probability is p is given by

$$R_n^*(c, p, r) = \max \left\{ rm : m \in \{0, 1, \dots, n\}, \sum_{\ell=1}^m W_{(\ell,n)} \leq c \text{ and } W_{(\ell,n)} \in \mathcal{W} \text{ for all } \ell \in [n] \right\}. \quad (2.7)$$

Earlier work has considered unitary rewards and deterministic arrivals by studying the random variable $R_n^*(c, 1, 1)$. First along this line of research, Coffman et al. (1987) showed that

$$R_n^*(c, 1, 1) \sim nF(\epsilon_n(c)) \text{ in probability as } n \rightarrow \infty,$$

provided that the weight distribution F is continuous, strictly increasing in w when $F(w) < 1$, and $F(w) \sim Aw^\alpha$ as $w \rightarrow 0$ for some $A, \alpha > 0$. Four years later, Bruss and Robertson (1991) proved that the same result holds under more general conditions, and Boshuizen and Kertz (1999) established the asymptotic normality of $R_n^*(c, 1, 1)$ after the usual centering and scaling for different classes of

weight distribution F . Lemma 4.1 in Bruss and Robertson (1991) is particularly relevant to our discussion here since it tells us that

$$\mathbb{E}[R_n^*(c, 1, 1)] \leq nF(\epsilon_n(c)) \quad \text{for all } n \geq 1.$$

Here, we generalize this result by accounting for Bernoulli arrivals with probability $p \in (0, 1]$ and rewards equal to $r > 0$. Specifically, we show that

$$\mathbb{E}[R_n^*(c, p, r)] \leq nprF(\epsilon_{np}(c)) \quad \text{for all } n \geq 1.$$

Our proof relies on the observation that the offline-sort algorithm (2.7) can be equivalently described as an algorithm that selects items with weight that is below some threshold. For any given realization W_1, W_2, \dots, W_n , the offline-sort algorithm selects $N_n^* \equiv R_n^*(c, p, 1)$ items so one can compute the value $W_{(N_n^*, n)}$ of the largest weight that is selected for inclusion, and one can then select all of the items $i \in [n]$ that have weight $W_i \leq W_{(N_n^*, n)}$. A shortcoming of this interpretation is that one needs to know the realization of the weight W_i (as well as the realizations of all of the other weights) to compute the threshold $W_{(N_n^*, n)}$. As it turns out, this is not needed in general. The next lemma shows that there is a thresholding algorithm that makes the same selections of offline sort, but in which the threshold used to decide whether to select an item is computed without using the information about that item's weight.

Lemma 2.1 (Threshold policy equivalence). *Let $W_{(1,n)} \leq W_{(2,n)} \leq \dots \leq W_{(n,n)}$ be the order statistics of $\mathcal{W} \equiv \{W_1, W_2, \dots, W_n\}$ and, for $i \in [n]$, let $W_{(1,n-1)} \leq W_{(2,n-1)} \leq \dots \leq W_{(n-1,n-1)}$ be the order statistics of $\mathcal{W}_i = \mathcal{W} \setminus \{W_i\}$. Then, for*

$$\tau_{n-1}^i = \max \left\{ m \in \{0, 1, \dots, n-1\} : \sum_{\ell=1}^m W_{(\ell, n-1)} \leq c \text{ and } W_{(\ell, n-1)} \in \mathcal{W}_i \text{ for all } \ell \in [n-1] \right\} \quad (2.8)$$

and $N_n^* \equiv R_n^*(c, p, 1)$, we have that

$$W_i \leq W_{(N_n^*, n)} \quad \text{if and only if} \quad W_i \leq h(\mathcal{W}_i) \equiv \max \left\{ W_{(\tau_{n-1}^i, n-1)}, c - \sum_{\ell=1}^{\tau_{n-1}^i} W_{(\ell, n-1)} \right\}. \quad (2.9)$$

In turn, it follows that

$$R_n^*(c, p, r) = \sum_{i=1}^n r \mathbb{1} \{W_i \leq h(\mathcal{W}_i)\}. \quad (2.10)$$

Proof. The equivalence (2.10) is an obvious consequence of (2.9), so we focus on proving the latter.

If $N_n^* = n$ we have that $\tau_{n-1}^i = n - 1$ and $W_i \leq c - \sum_{\ell=1}^{\tau_{n-1}^i} W_{(\ell, n-1)}$ for all $i \in [n]$, so equivalence (2.9) immediately follows. Instead, if $N_n^* < n$ the proof of (2.9) requires more work. As a warm-up we note that since the sets \mathcal{W} and \mathcal{W}_i differ only in one element, then

$$W_{(\ell, n)} \leq W_{(\ell, n-1)} \leq W_{(\ell+1, n)} \quad \text{for all } \ell \in [n-1]. \quad (2.11)$$

If we now recall the definitions of τ_{n-1}^i and N_n^* and use the inequalities above we obtain that

$$\sum_{\ell=1}^{\tau_{n-1}^i} W_{(\ell, n)} \leq \sum_{\ell=1}^{\tau_{n-1}^i} W_{(\ell, n-1)} \leq c \quad \text{and} \quad \sum_{\ell=1}^{N_n^*-1} W_{(\ell, n-1)} \leq \sum_{\ell=1}^{N_n^*-1} W_{(\ell+1, n)} \leq \sum_{\ell=1}^{N_n^*} W_{(\ell, n)} \leq c.$$

These two bounds respectively tell us that the offline-sort algorithm on \mathcal{W} selects at least τ_{n-1}^i observations, and that the same algorithm on \mathcal{W}_i selects at least $N_n^* - 1$ items. Thus, it follows that

$$N_n^* - 1 \leq \tau_{n-1}^i \leq N_n^*,$$

and we use these bounds to prove the equivalence (2.9).

If. We now suppose that $W_i \leq h(\mathcal{W}_i) \equiv \max \{W_{(\tau_{n-1}^i, n-1)}, c - \sum_{\ell=1}^{\tau_{n-1}^i} W_{(\ell, n-1)}\}$, and we seek to show that $W_i \leq W_{(N_n^*, n)}$. We consider two cases, one per each possible realization of τ_{n-1}^i .

Case 1: $\tau_{n-1}^i = N_n^* - 1$. If $\tau_{n-1}^i = N_n^* - 1$ then the definition of τ_{n-1}^i in (2.8) tells us that

$$c - \sum_{\ell=1}^{N_n^*-1} W_{(\ell, n-1)} < W_{(N_n^*, n-1)},$$

so if we apply the right inequality of (2.11) to $\ell = N_n^* - 1$ and $\ell = N_n^*$, we obtain that

$$W_{(N_n^*-1, n-1)} \leq W_{(N_n^*, n)} \quad \text{and} \quad c - \sum_{\ell=1}^{N_n^*-1} W_{(\ell, n-1)} < W_{(N_n^*+1, n)}. \quad (2.12)$$

If $W_{(N_n^*, n)} = W_{(N_n^*+1, n)}$ then the two inequalities in (2.12) give us that

$$h(\mathcal{W}_i) = \max \{W_{(N_n^*-1, n-1)}, c - \sum_{\ell=1}^{N_n^*-1} W_{(\ell, n-1)}\} \leq W_{(N_n^*, n)},$$

so we also have that $W_i \leq W_{(N_n^*, n)}$. On the other hand, if $W_{(N_n^*, n)} < W_{(N_n^*+1, n)}$ then the bounds in (2.12) imply that $h(\mathcal{W}_i) < W_{(N_n^*+1, n)}$, so we obtain from $W_i \leq h(\mathcal{W}_i)$ that $W_i \leq W_{(N_n^*, n)}$.

Case 2: $\tau_{n-1}^i = N_n^*$. The left inequality of (2.11) with $\ell = N_n^*$ tells us that we have two sub-cases to consider here: (i) when $W_{(N_n^*, n)}$ is equal to $W_{(N_n^*, n-1)}$, and (ii) when $W_{(N_n^*, n)}$ is strictly smaller than $W_{(N_n^*, n-1)}$. In the first sub-case, if $\tau_{n-1}^i = N_n^*$ and $W_{(N_n^*, n)} = W_{(N_n^*, n-1)}$, then the first N_n^* order statistics of \mathcal{W} and of \mathcal{W}_i agree and $c - \sum_{\ell=1}^{N_n^*} W_{(\ell, n-1)} = c - \sum_{\ell=1}^{N_n^*} W_{(\ell, n)} < W_{(N_n^*+1, n)}$. Thus, if $W_{(N_n^*, n)} = W_{(N_n^*+1, n)}$ then $h(\mathcal{W}_i) = \max\{W_{(N_n^*, n)}, c - \sum_{\ell=1}^{N_n^*} W_{(\ell, n)}\} = W_{(N_n^*, n)}$, and we are done. Otherwise, if $W_{(N_n^*, n)} < W_{(N_n^*+1, n)}$ then $h(\mathcal{W}_i) < W_{(N_n^*+1, n)}$ so that $W_i \leq h(\mathcal{W}_i) < W_{(N_n^*+1, n)}$ implies that $W_i \leq W_{(N_n^*, n)}$. In the second sub-case, if $\tau_{n-1}^i = N_n^*$ and $W_{(N_n^*, n)} < W_{(N_n^*, n-1)}$ then we have that $W_i = W_{(N_n^*, n)}$, and the result follows.

Only If. We now suppose that $W_i \leq W_{(N_n^*, n)}$, and we show that

$$W_i \leq h(\mathcal{W}_i) \equiv \max\left\{W_{(\tau_{n-1}^i, n-1)}, c - \sum_{\ell=1}^{\tau_{n-1}^i} W_{(\ell, n-1)}\right\}$$

by proving that $W_{(N_n^*, n)} \leq h(\mathcal{W}_i)$. Just as before, we consider separately the two possible realizations of τ_{n-1}^i .

Case 1: $\tau_{n-1}^i = N_n^* - 1$. We have two sub-cases to consider here. First, if $W_{(N_n^*, n)} \leq W_{(N_n^*-1, n-1)}$ then the lower bound $W_{(N_n^*, n)} \leq h(\mathcal{W}_i)$ is trivial. Second, if $W_{(N_n^*-1, n-1)} < W_{(N_n^*, n)}$ we show that the right maximand is bounded below by $W_{(N_n^*, n)}$. In this instance, the first $N_n^* - 1$ order statistic of \mathcal{W} and \mathcal{W}_i agree so the definition of N_n^* gives us that $W_{(N_n^*, n)} \leq c - \sum_{\ell=1}^{N_n^*-1} W_{(\ell, n)} = c - \sum_{\ell=1}^{N_n^*-1} W_{(\ell, n-1)}$, and we are done.

Case 2: $\tau_{n-1}^i = N_n^*$. If $\tau_{n-1}^i = N_n^*$ the left inequality of (2.11) tells us that $W_{(N_n^*, n)} \leq W_{(N_n^*, n-1)}$, so the lower bound $W_{(N_n^*, n)} \leq h(\mathcal{W}_i)$ immediately follows.

□

The representation (2.10) for $R_n^*(c, p, r)$ provides us with an easy way for proving the upper bound $\mathbb{E}[R_n^*(c, p, r)] \leq nprF(\epsilon_{np}(c))$. We just need to note that the expected total reward collected by the offline-sort algorithm is bounded above by the solution of some appropriate optimization problem. Our argument does not require independence of item weights. The threshold equivalence of Lemma 2.1 holds on every sample path, and the relaxation that follows only uses properties of the weight distribution F and of the arrival probability p (see also Steele, 2016).

Proposition 2.1 (Prophet upper bound). *Consider a knapsack problem with capacity $0 \leq c < \infty$ and with items that arrive over $1 \leq n < \infty$ periods according to a Bernoulli process with arrival probability $p \in (0, 1]$. If the items have rewards equal to r and weights with continuous distribution F , then for $\epsilon_{np}(c) = \sup \left\{ \epsilon \in [0, \infty) : \int_0^\epsilon w dF(w) \leq \frac{c}{np} \right\}$ we have that*

$$\mathbb{E}[R_n^*(c, p, r)] \leq nprF(\epsilon_{np}(c)). \quad (2.13)$$

Proof. To prove inequality (2.13), we begin with two easy cases. If $c = 0$ then $R_n^*(0, p, r) = 0$, and the bound (2.13) is trivial. Similarly, if $\mu = \mathbb{E}[W_1 \mid B_1 = 1] = \int_0^\infty w dF(w)$ and $np\mu < c < \infty$ then the definition of the function $\epsilon_{np}(c)$ tells us that $\epsilon_{np}(c) = +\infty$ so $F(\epsilon_{np}(c)) = 1$ and the bound (2.13) is again trivial because $R_n^*(c, p, r) \leq \sum_{i=1}^n r \mathbb{1}\{W_i < \infty\}$ for all $c \in [0, \infty)$, and this last right-hand side has expected value equal to npr .

Next, we consider the case in which $0 < c \leq np\mu$. If $\mathcal{W}_i \equiv \{W_1, \dots, W_{i-1}, W_{i+1}, \dots, W_n\}$ and $\mathcal{G}_i = \sigma\{\mathcal{W}_i\}$ is the σ -field generated by the sample \mathcal{W}_i , then we obtain from Lemma 2.1 and from the definition (2.7) that for each $i \in [n]$ there is a \mathcal{G}_i -measurable threshold $h(\mathcal{W}_i)$ such that one has the representation as well as the capacity constraint

$$R_n^*(c, p, r) = \sum_{i=1}^n r \mathbb{1}\{W_i \leq h(\mathcal{W}_i)\} \quad \text{and} \quad \sum_{i=1}^n W_i \mathbb{1}\{W_i \leq h(\mathcal{W}_i)\} \leq c.$$

In turn, we can obtain an upper bound for $\mathbb{E}[R_n^*(c, p, r)]$ by maximizing the sum $\sum_{i=1}^n \mathbb{E}[r \mathbb{1}\{W_i \leq h_i\}]$ over all thresholds (h_1, h_2, \dots, h_n) that satisfy an analogous capacity constraint and that have the same measurability property. Formally, we have the inequality

$$\begin{aligned} \mathbb{E}[R_n^*(c, p, r)] &\leq \max_{(h_1, \dots, h_n)} \sum_{i=1}^n \mathbb{E}[r \mathbb{1}\{W_i \leq h_i\}] \\ \text{s.t.} \quad &\sum_{i=1}^n W_i \mathbb{1}\{W_i \leq h_i\} \leq c \quad \text{almost surely} \\ &h_i \in \mathcal{G}_i \quad \text{for all } i \in [n]. \end{aligned} \quad (2.14)$$

Since $\epsilon_{np}(c) > 0$ and because the capacity constraint holds almost surely (and thus also in expecta-

tion), we have the further upper bound

$$\begin{aligned} \mathbb{E}[R_n^*(c, p, r)] &\leq \max_{(h_1, \dots, h_n)} \sum_{i=1}^n \mathbb{E}[r \mathbb{1}\{W_i \leq h_i\} \{1 - \epsilon_{np}^{-1}(c) W_i\}] + cr \epsilon_{np}^{-1}(c) \\ \text{s.t.} \quad &\sum_{i=1}^n W_i \mathbb{1}\{W_i \leq h_i\} \leq c \quad \text{almost surely} \\ &h_i \in \mathcal{G}_i \quad \text{for all } i \in [n]. \end{aligned} \quad (2.15)$$

Because h_i is \mathcal{G}_i -measurable, an application of the tower property gives us that

$$\mathbb{E}[\mathbb{E}[r \mathbb{1}\{W_i \leq h_i\} \{1 - \epsilon_{np}^{-1}(c) W_i\} \mid \mathcal{G}_i]] = pr \mathbb{E}\left[\int_0^{h_i} \{1 - \epsilon_{np}^{-1}(c) w\} dF(w)\right],$$

so, after we drop the two constraints in (2.15) we obtain that

$$\mathbb{E}[R_n^*(c, p, r)] \leq \mathfrak{p}^* = \max_{(h_1, \dots, h_n)} \sum_{i=1}^n pr \mathbb{E}\left[\int_0^{h_i} \{1 - \epsilon_{np}^{-1}(c) w\} dF(w)\right] + cr \epsilon_{np}^{-1}(c). \quad (2.16)$$

The maximization problem on the right-hand side is separable, and it is simple task to verify that the quantity $\mathbb{E}\left[\int_0^{h_i} \{1 - \epsilon_{np}^{-1}(c) w\} dF(w)\right]$ is maximized by setting $h_i = \epsilon_{np}(c)$ almost surely and for all $i \in [n]$. Thus, it follows that

$$\begin{aligned} \mathfrak{p}^* &= \sum_{i=1}^n \max_{h_i} pr \mathbb{E}\left[\int_0^{h_i} \{1 - \epsilon_{np}^{-1}(c) w\} dF(w)\right] + cr \epsilon_{np}^{-1}(c) \\ &= npr \left\{ F(\epsilon_{np}(c)) - \epsilon_{np}^{-1}(c) \left[\int_0^{\epsilon_{np}(c)} w dF(w) - \frac{c}{np} \right] \right\}. \end{aligned}$$

The integral representation (2.3) then tells us that the second summand is equal to zero, so after we recall (2.16) we obtain that

$$\mathbb{E}[R_n^*(c, p, r)] \leq \mathfrak{p}^* = npr F(\epsilon_{np}(c)) \quad \text{for all } 0 < c \leq np\mu,$$

completing the proof of (2.13). \square

2.3 The Reoptimized Policy $\hat{\pi}$ and Its Value Function

In the course of proving Proposition 2.1, we observed that if $\mathcal{G}_i = \sigma\{W_1, \dots, W_{i-1}, W_{i+1}, \dots, W_k\}$ is the σ -field generated by the sample $\{W_1, \dots, W_{i-1}, W_{i+1}, \dots, W_k\}$, then the expected value of the

offline solution $R_k^*(x, p, r)$ satisfies the upper bound

$$\begin{aligned} \mathbb{E}[R_k^*(x, p, r)] &\leq \max_{(h_1, \dots, h_k)} \sum_{i=1}^k \mathbb{E}[r \mathbb{1}\{W_i \leq h_i\}] \\ \text{s.t.} \quad &\sum_{i=1}^k W_i \mathbb{1}\{W_i \leq h_i\} \leq x \quad \text{almost surely} \\ &h_i \in \mathcal{G}_i \quad \text{for all } i \in [k]. \end{aligned}$$

We also noticed that the optimization problem on the right-hand side can be relaxed by first adding to its objective the quantity $\epsilon_{kp}^{-1}(x)r \left\{ x - \mathbb{E} \left[\sum_{i=1}^k W_i \mathbb{1}\{W_i \leq h_i\} \right] \right\} \geq 0$, and then by dropping the two constraints. This then gives us the further upper bound

$$\mathbb{E}[R_k^*(x, p, r)] \leq \max_{(h_1, \dots, h_k)} \sum_{i=1}^k pr \mathbb{E} \left[\int_0^{h_i} \{1 - \epsilon_{kp}^{-1}(c)w\} dF(w) \right] + xr \epsilon_{kp}^{-1}(x), \quad (2.17)$$

which is maximized by setting $h_i = \epsilon_{kp}(x)$ for all $i \in [k]$. We can now use this reoptimized solution for all $x \in [0, \infty)$ and all $1 \leq k < \infty$ to construct the online feasible threshold policy $\hat{\pi} \in \Pi(n, c, p)$. Specifically, since $\epsilon_{kp}(x)$ may exceed x , we set for $p \in (0, 1]$

$$\hat{h}_k(x) = \min\{x, \epsilon_{kp}(x)\}, \quad (2.18)$$

and we define the reoptimized policy $\hat{\pi}$ through the threshold $\{\hat{h}_n, \hat{h}_{n-1}, \dots, \hat{h}_1\}$. Thus, if the remaining capacity is x when item i is first presented, then item i is selected if and only if its weight $W_i \leq \hat{h}_{n-i+1}(x)$.

In turn, the threshold functions $\{\hat{h}_k : 1 \leq k < \infty\}$ induce a sequence of value functions $\{\hat{v}_k : [0, \infty) \rightarrow \mathbb{R}_+ : 0 \leq k < \infty\}$ such that $\hat{v}_k(x)$ represents the expected reward to-go of the reoptimized policy when there are k remaining periods and the current level of remaining knapsack capacity is x . If $\hat{v}_0(x) = 0$ for all $x \in [0, \infty)$, then the value $\hat{v}_k(x)$ is given by the recursion

$$\begin{aligned} \hat{v}_k(x) &= p \left(1 - F(\hat{h}_k(x)) \right) \hat{v}_{k-1}(x) + p \int_0^{\hat{h}_k(x)} \{r + \hat{v}_{k-1}(x - w)\} dF(w) + (1 - p) \hat{v}_{k-1}(x) \\ &= \left(1 - pF(\hat{h}_k(x)) \right) \hat{v}_{k-1}(x) + p \int_0^{\hat{h}_k(x)} \{r + \hat{v}_{k-1}(x - w)\} dF(w). \end{aligned} \quad (2.19)$$

By setting the number of remaining periods to n and the knapsack capacity to c , we find that

$$\hat{v}_n(c) = \mathbb{E} \left[R_n^{\hat{\pi}}(c, p, r) \right].$$

To verify the validity of the recursion (2.19), we condition on what happens in the k th-to-last period. With probability $1 - p$ the arriving item has arbitrarily large weight (equivalently, no item arrives), the number of the remaining periods decreases to $k - 1$ and the level of remaining capacity, x , stays the same. This then yields the term $(1 - p)\hat{v}_{k-1}(x)$ in the first line of (2.19). On the other hand, with probability p the arriving item has weight distribution F , and we can further condition on its realization, w . If $w > \hat{h}_k(x)$ then the item is rejected, the level of remaining capacity does not change, and the number of remaining periods decreases by one. That is, if the item is rejected, the expected reward to-go is given by $\hat{v}_{k-1}(x)$ and, since rejections happen with probability $p(1 - F(\hat{h}_k(x)))$, we recover the first summand on the top line of (2.19). On the other hand, if $w \leq \hat{h}_k(x)$ the k th-to-last item is included in the knapsack. Such a decision produces an immediate reward of r , and it depletes w units of capacity. The new remaining capacity then becomes $x - w$, and the number of remaining periods decreases to $k - 1$. The decision maker's payoff for including this item is then given by $r + \hat{v}_{k-1}(x - w)$ and, by integrating this payoff against the measure $p dF(w)$ for $w \in [0, \hat{h}_k(x)]$, we find the second summand on the first line of the recursion (2.19).

The reoptimized heuristic $\hat{\pi}$ then takes the solution of the offline relaxation (2.17) and turns it into an online algorithm through the threshold \hat{h}_k given in (2.18). This direct link provides us with enough tractability to be able to quantify the difference in expected performance between the reoptimized heuristic and the offline solution and—as a result—to prove the logarithmic regret bound. Instead, the optimal dynamic programming policy cannot be expressed explicitly and it lacks of the regularity needed to make any meaningful analytical progress. However, we note here that both the reoptimized heuristic and the optimal dynamic programming policy can be computed numerically in polynomial time, and we refer the reader to Section 2.6 for more details on our numerical work.

2.4 On the Typical Class

The weight distribution F plays a crucial role in the study of the performance of optimal and near-optimal item selections for the dynamic and stochastic knapsack problem with equal rewards. Because the weights are not equal, the remaining capacity process exhibits substantial randomness, and this may lead to unexpected behavior. As such, regularity conditions on the weight distribution

F are commonplace in the related literature. For instance, Coffman et al. (1987) only consider distributions F such that $F(w) \sim Aw^\alpha$ as $w \rightarrow 0$ for some $A, \alpha > 0$, while Bruss and Robertson (1991) expand this class to include all of the weight distributions F such that $\limsup_{w \rightarrow 0^+} F(\lambda w)/F(w) < 1$. Furthermore, Papastavrou et al. (1996, Section 5) show that one must require concavity of F to obtain structural properties such as monotonicity of the optimal threshold functions and concavity of the optimal value functions.

Here, we consider distributions that belong to the typical class characterized in Definition 2.1. As we mentioned earlier, this class is broad enough to include most well-known non-negative continuous distributions. Such breadth comes from the fact that Conditions (2.4) and (2.5) in Definition 2.1 must hold only on $(0, \bar{w})$ for some $\bar{w} > 0$, and that one has the flexibility of choosing different parameter \bar{w} for different distribution F . For instance, the uniform distribution $f(w) = \mathbb{1}\{w \in (0, 1)\}$ and the exponential distribution $f(w) = \alpha e^{-\alpha w} \mathbb{1}\{w > 0\}$ are both typical, but they require different choices of \bar{w} . For the uniform distribution, Conditions (2.4) and (2.5) hold on all of its support and one can choose $\bar{w} = 1$, while for the exponential distribution, Condition (2.5) holds only on $(0, 3/\alpha)$ and one can set $\bar{w} = 3/\alpha$. Similarly, one can check that the truncated normal distribution on $(0, b)$ with density $f(w) = A \exp\{-(w - v)^2/(2\varsigma^2)\} \mathbb{1}\{w \in (0, b)\}$ for $v \in \mathbb{R}$, $\varsigma > 0$, and A being the appropriate normalizing constant, is typical with $\bar{w} = \min\{\frac{1}{2}(v + \sqrt{v^2 + 12\varsigma^2}), b\}$. The truncated logistic distribution on $(0, b)$ and the logit-normal distribution are additional examples of typical distributions, though the respective \bar{w} 's have to do with the smallest positive root of related transcendental equations. The families of distributions listed below also belong to the typical class.

1. *Power distributions.* Distributions such that $F(w) = Aw^\alpha$ for some $A, \alpha > 0$ on $(0, \bar{w})$ are typical. Condition (2.4) is immediately verified. The function $w^3 f(w) = A\alpha w^{\alpha+2}$ is increasing because $A, \alpha > 0$, so (2.5) holds as well.
2. *Convex distributions.* Distributions F that are convex in a neighborhood of 0 and that have continuous density f are typical. Convexity tells us that $F(\lambda w) \leq F(w)\lambda$ so (2.4) follows. Furthermore, convexity also gives us that the density f is non-decreasing, so (2.5) is verified.
3. *Mixtures of typical distributions.* The class of typical distributions is closed under mixture. If F and G are two typical distributions and $\beta \in [0, 1]$ then it is easy to see that the mixture distribution $\beta F + (1 - \beta)G$ is also typical.

It is important to note, however, that one can construct examples of distributions that do not belong to the typical class. For instance, the distribution $F(w) = \frac{\log \bar{w}}{\log w}$ for $\bar{w} < 1$ and $w \in (0, \bar{w})$ is an example that satisfies Condition (2.5) but violates Condition (2.4). For a fixed $0 < \lambda < 1$, one can easily check that

$$\limsup_{w \rightarrow 0^+} \frac{F(\lambda w)}{F(w)} = \limsup_{w \rightarrow 0^+} \frac{\log w}{\log \lambda + \log w} = 1,$$

so Condition (2.4) fails to hold. On the other hand, the function $w^3 f(w) = -\frac{w^2 \log \bar{w}}{(\log w)^2}$ is increasing on $(0, \bar{w})$ and Condition (2.5) is satisfied.

The distribution $F(w) = A \int_0^w \{\sin(1/u)\}^2 du$ for $w \in (0, \bar{w})$ and $A = (\int_0^{\bar{w}} \{\sin(1/u)\}^2 du)^{-1} > 0$ is an example that satisfies Condition (2.4) while violating Condition (2.5). In fact, one has that the limit

$$\limsup_{w \rightarrow 0^+} \frac{F(\lambda w)}{F(w)} = \lambda < 1,$$

but the function $w^3 f(w) = Aw^3 \{\sin(1/w)\}^2$ oscillates infinitely many times in a (positive) neighborhood of zero, so the monotonicity (2.5) fails to hold.

We conclude this section by observing that Condition (2.4) regarding the behavior of F at zero is equivalent to the condition required by Bruss and Robertson (1991), and by proving that we can equivalently state it as a property of the ratio $wF(w)/\int_0^w u dF(u)$. This equivalent property will be important to our analysis.

Lemma 2.2 (Equivalence of CDF Conditions). *There are constants $0 < \lambda < 1$ and $0 < \gamma < 1$ and a value $\bar{w} > 0$ such that*

$$\frac{F(\lambda w)}{F(w)} \leq \gamma < 1 \quad \text{for all } w \in (0, \bar{w}) \quad (2.20)$$

if and only if there is a constant $1 < M < \infty$ such that

$$\frac{wF(w)}{\int_0^w u dF(u)} \leq M < \infty \quad \text{for all } w \in (0, \bar{w}). \quad (2.21)$$

Proof. If. Suppose there is a constant $1 < M < \infty$ such that condition (2.21) holds. Next, note that for any $\lambda \in (0, 1)$ and any $w \in (0, \bar{w})$ one has the bounds

$$0 \leq \int_0^w u dF(u) \leq \lambda w \int_0^{\lambda w} dF(u) + w \int_{\lambda w}^w dF(u) = wF(w) - wF(\lambda w)(1 - \lambda),$$

so it follows that

$$\frac{F(w)}{F(w) - F(\lambda w)(1 - \lambda)} \leq \frac{wF(w)}{\int_0^w u dF(u)}.$$

In turn, condition (2.21) tells us that there is $1 < M < \infty$ such that the right-hand side above is bounded by M so, after rearranging, we obtain that

$$\frac{F(\lambda w)}{F(w)} \leq \frac{M - 1}{M(1 - \lambda)} \quad \text{for all } w \in (0, \bar{w}).$$

Condition (2.20) then follows after one chooses any $\lambda < M^{-1}$ and sets $\gamma = (M - 1)/[M(1 - \lambda)] < 1$.

Only if. Suppose that there are constants $0 < \lambda < 1$ and $0 < \gamma < 1$ such that condition (2.20) holds for some $\bar{w} > 0$. Then we have that

$$0 < 1 - \gamma \leq 1 - \frac{F(\lambda w)}{F(w)} = \int_{\lambda w}^w \frac{dF(u)}{F(w)} \quad \text{for all } w \in (0, \bar{w}).$$

Moreover, if we multiply both sides by λw and use the fact that $\lambda w \leq u$ for all $u \in (\lambda w, w)$ we also have that

$$\lambda w(1 - \gamma) \leq \int_{\lambda w}^w \frac{\lambda w}{F(w)} dF(u) \leq \int_0^w \frac{u}{F(w)} dF(u).$$

Next, we divide both sides by w and rearrange to obtain that

$$\frac{wF(w)}{\int_0^w u dF(u)} \leq \frac{1}{\lambda(1 - \gamma)} \quad \text{for all } w \in (0, \bar{w}),$$

so condition (2.21) follows by setting $M = [\lambda(1 - \gamma)]^{-1}$, and the proof is now complete. \square

2.5 A Logarithmic Regret Bound

To prove that the regret grows at most logarithmically, we let

$$K = \left\lceil \frac{c}{p \int_0^{\bar{w}} w f(w) dw} \right\rceil \tag{2.22}$$

and focus on dynamic and stochastic knapsack problems with more than K periods. Of course, this is without loss of generality because the quantity K defined in (2.22) is a constant that does not depend on the number of periods n , so we can ignore the last K decisions without affecting our

regret bound. When $k \geq K$ we have (i) that $\epsilon_{kp}(x) \leq \bar{w}$ for all $x \in [0, c]$, and (ii) that the integral representation (2.3) always holds. Thus, we are focusing on problem instances in which we can use the properties of the typical class in full.

In our proof, we will repeatedly use the following two properties of the consumption function $\epsilon_{kp}(x)$. First, we obtain from definition (2.2) that the consumption functions are non-increasing in k . That is, for $p \in (0, 1]$ one has the monotonicity

$$\epsilon_{(k+1)p}(x) \leq \epsilon_{kp}(x) \quad \text{for all } x \in [0, \infty) \text{ and all } k \geq 1. \quad (2.23)$$

Second, provided that the weight distribution F has continuous density f , an application of the implicit function theorem gives us that the function $\epsilon_{kp}(x)$ is differentiable on $(0, kp\mu)$, and that its first derivative $\epsilon'_{kp}(x)$ is given by

$$\epsilon'_{kp}(x) = \frac{1}{kp\epsilon_{kp}(x)f(\epsilon_{kp}(x))} \quad \text{if } 0 < x < kp\mu. \quad (2.24)$$

The proof of the regret bound then comes in two parts. In the next section we derive several estimates that have to do with the weight distribution belonging to the typical class and with $k \geq K$, while in Section 2.5.2 we estimate the gap $kprF(\epsilon_{kp}(x)) - \hat{v}_k(x)$.

2.5.1 Preliminary Observations

When $k \geq K$ the properties that characterize typical weight distributions can be used to obtain general estimates that are crucial to our analysis. As a warm-up we obtain the following estimate on the mismatch between the probability of an item weight being smaller than the feasible threshold \hat{h}_k and the probability of the same weight being smaller than the consumption function ϵ_{kp} .

Lemma 2.3. *If the weight distribution F belongs to the typical class then there is $1 < M < \infty$ depending only on F such that*

$$\frac{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))}{x} \leq M \quad \text{for all } x \in (0, c], p \in (0, 1], \text{ and all } k \geq K \equiv \left\lceil \frac{c}{p \int_0^{\bar{w}} wf(w) dw} \right\rceil. \quad (2.25)$$

In turn, we also have that

$$F(\epsilon_{kp}(x)) - F(\hat{h}_k(x)) \leq \frac{M}{kp} \quad \text{for all } x \in [0, c], p \in (0, 1], \text{ and all } k \geq K. \quad (2.26)$$

Proof. The uniform bound (2.25) is essentially a restatement of inequality (2.21) in Lemma 2.2. If $x \in (0, c]$ and $k \geq K$, then we have that

$$\frac{x}{kp} \leq \frac{x}{Kp} \leq \frac{c}{Kp} \leq \int_0^{\bar{w}} wf(w) dw \leq \mu,$$

so the definition (2.2) of the consumption function $\epsilon_{kp}(x)$ and the equality (2.3) give us that

$$\epsilon_{kp}(x) \leq \bar{w} \quad \text{and} \quad \int_0^{\epsilon_{kp}(x)} wf(w) dw = \frac{x}{kp} \quad \text{for all } k \geq K \text{ and all } x \in (0, c]. \quad (2.27)$$

The two observations in (2.27) together with the bound (2.21) in which we replace w with $\epsilon_{kp}(x)$ then imply that

$$\frac{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))}{x} = \frac{\epsilon_{kp}(x)F(\epsilon_{kp}(x))}{\int_0^{\epsilon_{kp}(x)} uf(u) du} \leq M \quad \text{for all } k \geq K \text{ and } x \in (0, c],$$

concluding the proof of the uniform bound (2.25).

We now turn to inequality (2.26). If $x = 0$ then inequality (2.26) is obvious. Otherwise, if $x > 0$ we recall from (2.18) that $\hat{h}_k(x) = \min\{x, \epsilon_{kp}(x)\}$, so the left-hand side of (2.26) is equal to 0 when $\epsilon_{kp}(x) \leq x < \infty$, and inequality (2.26) is again trivial. Instead, if $0 < x < \epsilon_{kp}(x)$, we obtain from (2.25) that

$$F(\epsilon_{kp}(x)) \leq \frac{Mx}{kp\epsilon_{kp}(x)} \leq \frac{M}{kp} \quad \text{for all } k \geq K \text{ and } 0 < x < \epsilon_{kp}(x),$$

concluding the proof of the lemma. □

In the same spirit of Lemma 2.3, we can also estimate the difference in the probability of selecting an upcoming item as a function of the number of remaining periods.

Lemma 2.4. *For $p \in (0, 1]$, all $x \in [0, c]$, and all $k \geq K$ we have that*

$$F(\epsilon_{(k+1)p}(x)) - F(\epsilon_{kp}(x)) \leq -\frac{x}{k(k+1)p\epsilon_{kp}(x)}.$$

Proof. For $k \geq K$ the equality (2.3) and the monotonicity (2.23) give us the representation

$$\frac{x}{kp} - \frac{x}{(k+1)p} = \int_{\epsilon_{(k+1)p}(x)}^{\epsilon_{kp}(x)} wf(w) dw, \quad \text{for all } x \in [0, c].$$

If we now replace the integrand $wf(w)$ with the upper bound $\epsilon_{kp}(x)f(w)$ and rearrange, we obtain

$$\frac{x}{k(k+1)p} \leq \epsilon_{kp}(x) [F(\epsilon_{kp}(x)) - F(\epsilon_{(k+1)p}(x))],$$

completing the proof of the lemma. \square

Typical weight distributions are also nice because one can tightly approximate the difference $F(\epsilon_{kp}(x)) - F(\epsilon_{kp}(x-w))$ that accounts for the sensitivity in the remaining capacity of the probability of selecting the k th-to-last item. A formal estimate is given in the next proposition, and it constitutes a key step in our argument.

Proposition 2.2. *If $p \in (0, 1]$ and if the weight distribution F belongs to the typical class, then there is a constant $1 < M < \infty$ depending only on F such that one has the inequality*

$$1 - \frac{F(\epsilon_{kp}(x-w))}{F(\epsilon_{kp}(x))} \leq \frac{w^2}{x^2} (1 - M^{-1}) + \frac{w}{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))} \quad (2.28)$$

for all $w \in [0, x]$, $x \in (0, c]$, and all $k \geq K \equiv \left\lceil \frac{c}{p \int_0^w wf(w) dw} \right\rceil$.

The proof of Proposition 2.2 requires the following intermediate estimate.

Lemma 2.5 (Convexity upper bound). *If $p \in (0, 1]$ and if the weight distribution F has continuous density f then for all $k \geq K$, $x \in [0, c]$ and $y \in [0, 1]$ we have the integral representation*

$$kpF(\epsilon_{kp}(x)) - kpF(\epsilon_{kp}(x(1-y))) = \int_{x(1-y)}^x \frac{1}{\epsilon_{kp}(u)} du. \quad (2.29)$$

Moreover, if the distribution F belongs to the typical class the map $x \mapsto \epsilon_{kp}(x)^{-1}$ is convex on $(0, c)$, so we also have the upper bound

$$kpF(\epsilon_{kp}(x)) - kpF(\epsilon_{kp}(x(1-y))) \leq \frac{xy}{2} \left[\frac{1}{\epsilon_{kp}(x)} + \frac{1}{\epsilon_{kp}(x(1-y))} \right]. \quad (2.30)$$

Proof. Since the weight distribution F has continuous density and $\frac{c}{p\mu} \leq K \leq k$, we recall from (2.24) the first derivative

$$\epsilon'_{kp}(x) = \frac{1}{kp\epsilon_{kp}(x)f(\epsilon_{kp}(x))} \quad \text{for all } x \in (0, c).$$

The map $x \mapsto F(\epsilon_{kp}(x))$ is then differentiable on $(0, c)$, and one has that

$$(kpF(\epsilon_{kp}(x)))' = kp\epsilon'_{kp}(x)f(\epsilon_{kp}(x)) = \frac{1}{\epsilon_{kp}(x)} \quad \text{for all } x \in (0, c).$$

In turn, the fundamental theorem of calculus tells us that for $y \in [0, 1]$ we have the integral representation

$$kpF(\epsilon_{kp}(x)) - kpF(\epsilon_{kp}(x(1-y))) = \int_{x(1-y)}^x \frac{1}{\epsilon_{kp}(u)} du,$$

proving the first assertion of the lemma.

To check the convexity of the map $x \mapsto \epsilon_{kp}(x)^{-1}$, we use the expression of the first derivative (2.24) one more time to obtain for $k \geq K$ that

$$\left(\frac{1}{\epsilon_{kp}(x)} \right)' = -\frac{\epsilon'_{kp}(x)}{\epsilon_{kp}^2(x)} = -\frac{1}{kp\epsilon_{kp}^3(x)f(\epsilon_{kp}(x))}.$$

If F belongs to the typical class and $k \geq K$ then the monotonicity condition (2.5) implies that the first derivative $(1/\epsilon_{kp}(x))'$ is non-decreasing on $(0, c)$, so the map $x \mapsto \epsilon_{kp}(x)^{-1}$ is convex. This convexity property then provides us with a linear majorant

$$m_{kp}(u) = \frac{u-x}{yx} \left(\frac{1}{\epsilon_{kp}(x)} - \frac{1}{\epsilon_{kp}((1-y)x)} \right) + \frac{1}{\epsilon_{kp}(x)}$$

such that

$$\frac{1}{\epsilon_{kp}(u)} \leq m_{kp}(u) \quad \text{for all } u \in [(1-y)x, x].$$

The representation (2.29) and the integration of the majorant $m_{kp}(u)$ over $[(1-y)x, x]$ give us the upper bound (2.30), and the proof of the lemma follows. \square

We now have all of the estimates we need to complete the proof of Proposition 2.2.

Proof. (Proof of Proposition 2.2.) If $w = 0$ then inequality (2.28) is trivial. Otherwise, for $K \leq k < \infty$ we consider the function $g_k : (0, c] \times (0, 1] \rightarrow \mathbb{R}$ given by

$$g_k(x, y) = \frac{1}{y^2} \left\{ 1 - \frac{F(\epsilon_{kp}(x(1-y)))}{F(\epsilon_{kp}(x))} - \frac{xy}{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))} \right\},$$

and we note that inequality (2.28) follows by setting $y = w/x \leq 1$ and rearranging, provided that one has the uniform bound

$$g_k(x, y) \leq 1 - M^{-1} \quad \text{for all } x \in (0, c], y \in (0, 1], \text{ and } k \geq K. \quad (2.31)$$

The function $g_k(x, y)$ is differentiable with respect to y for any given $x \in (0, c]$, and the y -derivative of $g_k(x, y)$ can be written as

$$\frac{\partial}{\partial y} g_k(x, y) = \frac{2}{y^3 k p F(\epsilon_{kp}(x))} \left\{ \frac{xy}{2} \left[\frac{1}{\epsilon_{kp}(x)} + \frac{1}{\epsilon_{kp}(x(1-y))} \right] - k p F(\epsilon_{kp}(x)) + k p F(\epsilon_{kp}(x(1-y))) \right\}.$$

Since $\frac{2}{y^3 k p F(\epsilon_{kp}(x))} \geq 0$, inequality (2.30) of Lemma 2.5 then tells us that the y -derivative of $g_k(x, y)$ is non-negative so that the map $y \mapsto g_k(x, y)$ is non-decreasing in y for any given $x \in (0, c]$. In turn, we have that

$$g_k(x, y) \leq g_k(x, 1) = 1 - \frac{x}{k p \epsilon_{kp}(x) F(\epsilon_{kp}(x))},$$

so inequality (2.31) follows from the uniform bound (2.25), and the proof of the proposition is now complete. \square

2.5.2 Analysis of Residuals

To estimate the gap between the expected total reward collected by the reoptimized policy $\hat{\pi} \in \Pi(n, c, p)$ and the prophet upper bound $n p r F(\epsilon_{np}(c))$, we study appropriate residual functions. Specifically, we let

$$\rho_k(x) = k p r F(\epsilon_{kp}(x)) - \hat{v}_k(x) \quad \text{for } x \in [0, c] \text{ and } 1 \leq k \leq n \quad (2.32)$$

be the *residual function* when there are k remaining periods and the level of remaining capacity is x . The residual function $\rho_k(x)$ is continuous and defined on a compact interval, so if we maximize with respect to x we obtain the *maximal residual*

$$\bar{\rho}_k = \max_{0 \leq x \leq c} \rho_k(x) \quad \text{for } k \in [n]. \quad (2.33)$$

The second half of Theorem 2.1 is just a corollary of the following proposition, which verifies that the maximal residual $\bar{\rho}_n = O(\log n)$ as $n \rightarrow \infty$.

Proposition 2.3. *If the weight distribution F belongs to the typical class, then there is a constant $1 < M < \infty$ depending only on the distribution F , the arrival probability p , and the reward r such that the maximal residual*

$$\bar{\rho}_n = \max_{0 \leq x \leq c} \{nprF(\epsilon_{np}(x)) - \hat{v}_n(x)\} \leq M + M \log n \quad \text{for all } n \geq 1.$$

For the proof of this proposition we write the maximal residual $\bar{\rho}_n$ as a telescoping sum, and we obtain an appropriate upper bound for each summand. The upper bound follows from the following lemma.

Lemma 2.6. *If the weight distribution F belongs to the typical class, then there is a constant $1 < M < \infty$ that depends only on F and the reward r such that the difference*

$$\rho_{k+1}(x) - \bar{\rho}_k \leq \frac{M}{k+1} \quad \text{for all } x \in [0, c] \text{ and all } k \geq K.$$

Proof. The residual function $\rho_k(x)$ defined in (2.32) provides us with an alternative representation for the value function $\hat{v}_{k+1}(x)$ which gives us the expected total reward selected by policy $\hat{\pi}$ with $k+1$ periods remaining and current knapsack capacity x . Specifically, if we substitute $\hat{v}_k(x)$ with $kprF(\epsilon_{kp}(x)) - \rho_k(x)$ in the recursion (2.19), we then obtain that

$$\begin{aligned} \hat{v}_{k+1}(x) = & \{1 - pF(\hat{h}_{k+1}(x))\} \{kprF(\epsilon_{kp}(x)) - \rho_k(x)\} \\ & + p \int_0^{\hat{h}_{k+1}(x)} \{r + kprF(\epsilon_{kp}(x-w)) - \rho_k(x-w)\} f(w) dw. \end{aligned}$$

Next, if we replace the residuals $\rho_k(\cdot)$ with their maximal value $\bar{\rho}_k$ and rearrange, we obtain the lower bound

$$\begin{aligned} kprF(\epsilon_{kp}(x)) &+ prF(\hat{h}_{k+1}(x)) + kp^2r \int_0^{\hat{h}_{k+1}(x)} \{F(\epsilon_{kp}(x-w)) - F(\epsilon_{kp}(x))\} f(w) dw \quad (2.34) \\ &\leq \hat{v}_{k+1}(x) + \bar{\rho}_k. \end{aligned}$$

In turn, the definition (2.32) of the residual function tells us that

$$\rho_{k+1}(x) - \bar{\rho}_k = (k+1)prF(\epsilon_{(k+1)p}(x)) - (\hat{v}_{k+1}(x) + \bar{\rho}_k),$$

so if we replace the sum $\hat{v}_{k+1}(x) + \bar{\rho}_k$ with its lower bound (2.34) and rearrange, we obtain the

upper bound

$$\begin{aligned} \rho_{k+1}(x) - \bar{\rho}_k &\leq pr \left\{ (k+1)F(\epsilon_{(k+1)p}(x)) - kF(\epsilon_{kp}(x)) - F(\hat{h}_{k+1}(x)) \right\} \\ &\quad + kp^2 r F(\epsilon_{kp}(x)) \int_0^{\hat{h}_{k+1}(x)} \left\{ 1 - \frac{F(\epsilon_{kp}(x-w))}{F(\epsilon_{kp}(x))} \right\} f(w) dw. \end{aligned} \quad (2.35)$$

Next, we obtain from (2.28) that the integral that appears on the right-hand side of (2.35) satisfies the upper bound

$$\mathcal{I}_k(x) \equiv \int_0^{\hat{h}_{k+1}(x)} \left\{ 1 - \frac{F(\epsilon_{kp}(x-w))}{F(\epsilon_{kp}(x))} \right\} f(w) dw \leq \int_0^{\hat{h}_{k+1}(x)} \left\{ \frac{w^2}{x^2} (1 - M^{-1}) + \frac{w}{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))} \right\} f(w) dw.$$

For $w \in [0, \hat{h}_{k+1}(x)]$ we have the trivial bound $w^2 \leq w\hat{h}_{k+1}(x)$ so if we replace w^2 with its upper bound $w\hat{h}_{k+1}(x)$ on the right-hand side above and integrate we obtain that there is $1 < M < \infty$ depending only on F such that

$$\mathcal{I}_k(x) \leq \left[(1 - M^{-1}) \frac{\hat{h}_{k+1}(x)}{x^2} + \frac{1}{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))} \right] \int_0^{\hat{h}_{k+1}(x)} w f(w) dw.$$

We now multiply both sides by $kp^2 r F(\epsilon_{kp}(x))$ and simplify to obtain that

$$kp^2 r F(\epsilon_{kp}(x)) \mathcal{I}_k(x) \leq \left[kp^2 r F(\epsilon_{kp}(x)) (1 - M^{-1}) \frac{\hat{h}_{k+1}(x)}{x^2} + \frac{pr}{\epsilon_{kp}(x)} \right] \int_0^{\hat{h}_{k+1}(x)} w f(w) dw.$$

The definition of $\hat{h}_{k+1}(x) = \min\{x, \epsilon_{(k+1)p}(x)\}$ and the monotonicity (2.23) tell us that $\hat{h}_{k+1}(x) \leq \epsilon_{(k+1)p}(x) \leq \epsilon_{kp}(x)$, so we obtain a further upper bound if we replace the first $\hat{h}_{k+1}(x)$ on the last right-hand side with $\epsilon_{kp}(x)$ and the second one with $\epsilon_{(k+1)p}(x)$. When we perform these replacements and recall the equality (2.3), we find that

$$kp^2 r F(\epsilon_{kp}(x)) \mathcal{I}_k(x) \leq \frac{r(1 - M^{-1})}{k+1} \frac{kp\epsilon_{kp}(x)F(\epsilon_{kp}(x))}{x} + \frac{rx}{(k+1)\epsilon_{kp}(x)}.$$

If we now apply the uniform upper bound (2.25) to the first summand on the right-hand side, and rearrange, we obtain that

$$kp^2 r F(\epsilon_{kp}(x)) \mathcal{I}_k(x) \leq \frac{r(M-1)}{k+1} + \frac{rx}{(k+1)\epsilon_{kp}(x)}.$$

We now replace the last summand of (2.35) with the upper bound above and rearrange to obtain that

$$\begin{aligned} \rho_{k+1}(x) - \bar{\rho}_k &\leq \frac{r(M-1)}{k+1} + kpr \left\{ F(\epsilon_{(k+1)p}(x)) - F(\epsilon_{kp}(x)) + \frac{x}{k(k+1)p\epsilon_{kp}(x)} \right\} \\ &\quad + pr \left\{ F(\epsilon_{(k+1)p}(x)) - F(\hat{h}_{k+1}(x)) \right\}. \end{aligned}$$

Here, Lemma 2.4 tells us that the second summand on the right-hand side is non-positive, and inequality (2.26) tells us that there is $1 < M < \infty$ depending only on F such that the difference $F(\epsilon_{(k+1)p}(x)) - F(\hat{h}_{k+1}(x))$ is bounded above by $M/((k+1)p)$. When we assemble these observations, we finally find that

$$\rho_{k+1}(x) - \bar{\rho}_k \leq \frac{(2M-1)r}{k+1} \quad \text{for all } x \in [0, c] \text{ and all } k \geq K,$$

concluding the proof of the lemma. \square

We now have all of the tools we need to complete the proof of Proposition 2.3 that follows next.

Proof. (Proof of Proposition 2.3.) We write the maximal residual $\bar{\rho}_n$ in (2.33) as a telescoping sum and use the definition (2.32) of the residual function to obtain that

$$\bar{\rho}_n = \bar{\rho}_K + \sum_{k=K}^{n-1} \{\bar{\rho}_{k+1} - \bar{\rho}_k\} \leq K + \sum_{k=K}^{n-1} \{\bar{\rho}_{k+1} - \bar{\rho}_k\}.$$

Lemma 2.6 then tells us that

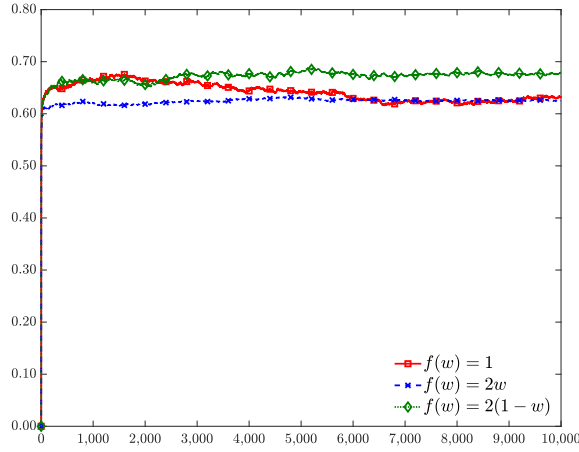
$$\bar{\rho}_{k+1} - \bar{\rho}_k \leq \frac{M}{k+1} \quad \text{for all } K \leq k \leq n,$$

so when we combine the last two observations we obtain that there is a constant $1 < M < \infty$ that depends only on F , p , and r such that

$$\bar{\rho}_n \leq M + M \log n,$$

just as needed. \square

Figure 2.1: Gap between the prophet upper bound and offline sort for three weight distributions.



Notes. Difference between the prophet upper bound, $nF(\epsilon_n(1))$, and the simulated average (with 100,000 trials) of the offline solution, $R_n^*(1, 1, 1)$, for three different distributions on the unit interval: $f(w) = \mathbb{1}\{w \in (0, 1)\}$, $f(w) = 2w\mathbb{1}\{w \in (0, 1)\}$ and $f(w) = 2(1-w)\mathbb{1}\{w \in (0, 1)\}$. In each case we take the arrival probability $p = 1$, the knapsack capacity $c = 1$, the reward $r = 1$, and we vary the number of periods $n \in \{1, 2, \dots, 10000\}$. The chart suggests that the gap between the prophet upper bound and the simulated average of the offline solution does not grow with n .

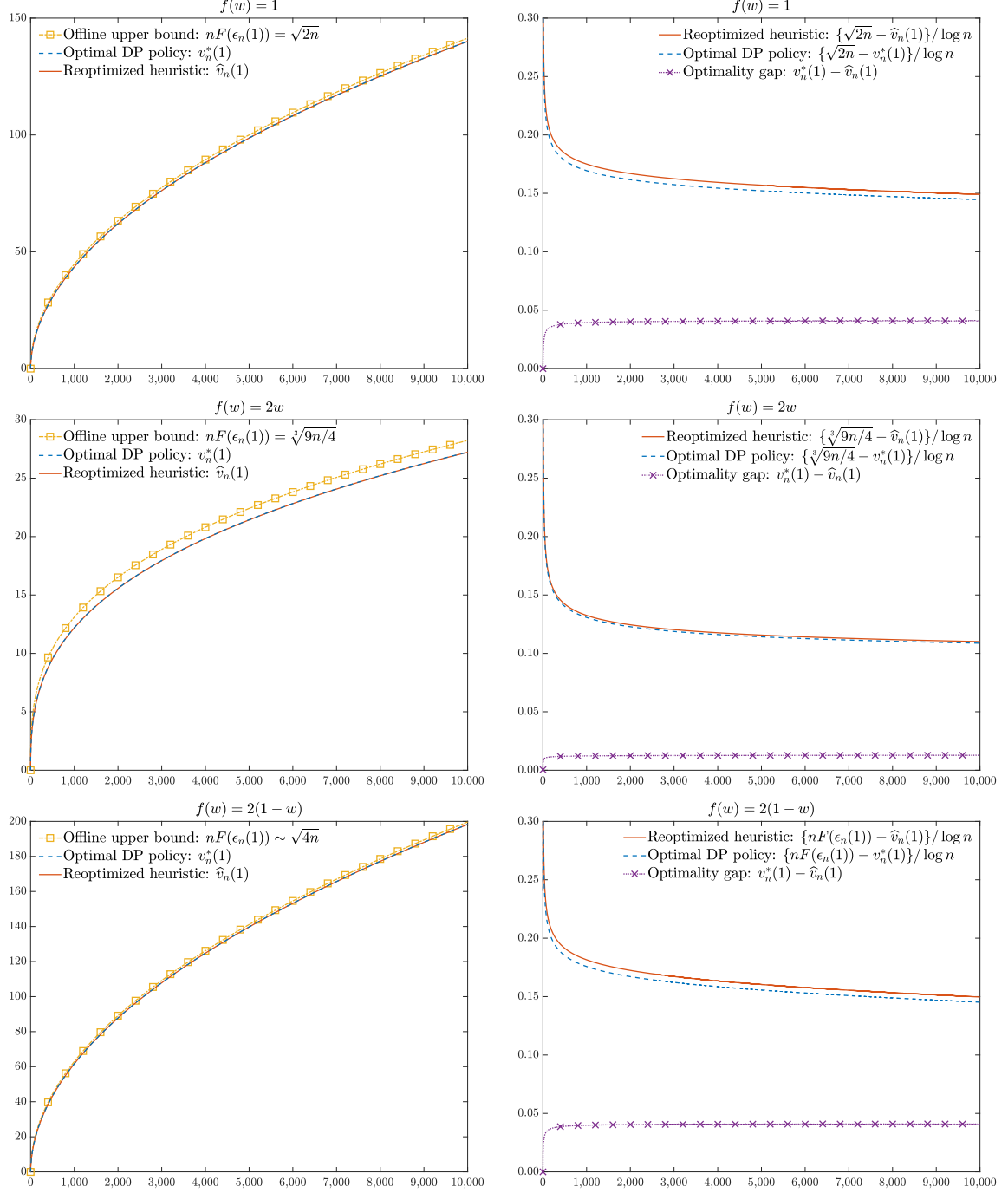
2.6 Numerical Experiments

Theorem 2.1 tells us that the regret of a dynamic and stochastic knapsack problem is at most logarithmic in n , provided that the weight distribution belongs to the typical class of Definition 2.1. While the actual order of the regret may—in principle—be smaller than what our bound predicts, we find numerically that this is not the case. In fact, we conjecture that the actual regret is $O(\log n)$ as $n \rightarrow \infty$ for most continuous weight distributions.

The work of Seksenbayev (2018) and Gnedin and Seksenbayev (2019) tells us that when the capacity, the reward, and the arrival probability are all equal to one, and the weight distribution is uniform on the unit interval, then the regret is asymptotic to $(\log n)/12$. In this section, we numerically investigate the actual order of the regret for two other weight distributions, while keeping the uniform as reference.

For our numerical examples, we solve the recursion (2.19) on a discretized state space with grid size 10^{-5} and obtain estimates for the reoptimized value function $\hat{v}_n(\cdot)$ for $n \in \{1, \dots, 10000\}$ and for different distributions F . On the same discretized state space and for the same weight distributions,

Figure 2.2: Value functions and scaled regret bounds for three weight distributions



Notes. The left plots display the prophet upper bound and the value functions of the optimal dynamic programming (DP) policy and of the reoptimized heuristic. The right plots show the regret bounds of the optimal policy and of the heuristic scaled by the logarithm of n , as well as the optimality gap. While the scaled regret bounds are bounded away from zero for large n , the optimality gap does not grow with n . Weights have densities on $(0, 1)$ respectively given by $f(w) = 1$, $f(w) = 2w$, and $f(w) = 2(1-w)$. Capacity $c = 1$, arrival probability $p = 1$, and reward $r = 1$. Discretized state space with grid size 10^{-5} .

we also solve numerically the Bellman recursion

$$\begin{aligned}
v_n^*(x) &= p(1 - F(x))v_{n-1}^*(x) + p \int_0^x \max\{r + v_{n-1}^*(x - w), v_{n-1}^*(x)\} f(w) dw + (1 - p)v_{n-1}^*(x) \\
&= (1 - pF(x))v_{n-1}^*(x) + p \int_0^x \max\{r + v_{n-1}^*(x - w), v_{n-1}^*(x)\} f(w) dw
\end{aligned} \tag{2.36}$$

with the initial condition $v_0^*(x) = 0$ for all $x \in [0, c]$, and we obtain estimates for the optimal value functions $v_n^*(\cdot)$ for $n \in \{1, \dots, 10000\}$. Finally, we simulate the average of the offline solution $R_n^*(c, p, r)$ and compare all of our numerical estimates with the prophet upper bound $nprF(\epsilon_{np}(c))$. Based on our numerical experiments, we observe that:

- (i) The gap $nprF(\epsilon_{np}(c)) - \mathbb{E}[R_n^*(c, p, r)]$ between the prophet upper bound and the offline solution is bounded by a constant that does not depend on n (see Figure 2.1);
- (ii) The regret bound $nprF(\epsilon_{np}(c)) - \hat{v}_n(c)$ for the reoptimized heuristic and the regret bound $nprF(\epsilon_{np}(c)) - v_n^*(c)$ for the optimal online policy grow logarithmically with n (Figure 2.2); and
- (iii) The optimality gap $v_n^*(c) - \hat{v}_n(c)$ is bounded by constant that is independent of n (Figure 2.2).

In turn, our numerical experiments suggests that the regrets (rather than the regret bounds) $\mathbb{E}[R_n^*(c, p, r)] - \hat{v}_n(c)$ and $\mathbb{E}[R_n^*(c, p, r)] - v_n^*(c)$ respectively of the reoptimized heuristic and of the optimal online policy are also logarithmic in n . In our numerical work, we consider instances of the dynamic and stochastic knapsack problem with reward $r = 1$, arrival probability $p = 1$, and capacity $c = 1$. We vary item weights by considering the three densities supported on the unit interval given by $f(w) = 1$, $f(w) = 2w$ and $f(w) = 2(1 - w)$ for $w \in (0, 1)$. The top left chart of Figure 2.2 plots the prophet upper bound $nF(\epsilon_n(1)) = \sqrt{2n}$ as well as the value function of the optimal policy, $v_n^*(1)$, and of the reoptimized heuristic, $\hat{v}_n(1)$, when the weight distribution is uniform on $(0, 1)$. Instead, the top right chart depicts the respective regret bounds scaled by the logarithm of n , as well as the optimality gap. In the chart we see that the scaled regret bounds (top two lines) are bounded away from zero for large n , implying that the regret bounds grow logarithmically. In contrast, the optimality gap (bottom line) appears not to grow with n .

The plots in the middle row of Figure 2.2 point to the same set of observations when the weights have density $f(w) = 2w \mathbb{1}\{w \in (0, 1)\}$ and the prophet upper bound is $nF(\epsilon_n(1)) = \sqrt[3]{9n/4}$.

Finally, the bottom two charts of Figure 2.2 consider item weights that have density $f(w) = 2(1 - w)\mathbb{1}\{w \in (0, 1)\}$. In this case, the prophet upper bound cannot be expressed in closed form, but one can show that $nF(\epsilon_n(1)) \sim \sqrt{4n}$ as $n \rightarrow \infty$. Nevertheless, also for this weight distribution the numerical analysis suggests that the regrets of the optimal policy and of the heuristic are both logarithmic in n , and that the optimality gap can be bounded by a constant independent of n .

2.7 On Weight Distributions with Multiple Types

In this section, we discuss how our logarithm regret bound generalizes to dynamic and stochastic knapsack problems with equal rewards and with independent random weights that belong to one of $J < \infty$ different types. We consider a multinomial arrival process with parameters $\mathbf{p} \equiv (p_0, p_1, \dots, p_J)$ where $p_j \in (0, 1]$ for all $j \in [J]$ and $p_0 = 1 - \sum_{i \in [J]} p_i \in [0, 1]$. Here, the parameter p_0 represents the probability of no item arriving (or, equivalently, the arrival probability of an item with arbitrarily large weight) and p_j , $j \in [J]$, is the arrival probability of an item with weight distribution F_j .

Upon arrival of an item the decision maker may see the type of the item or not. If the item types are *not* released, then she only sees the arriving weights that (conditional on an arrival occurring) are drawn from the mixture distribution

$$\tilde{F}(w) = \frac{1}{1 - p_0} \sum_{j \in [J]} p_j F_j(w) \quad \text{for all } w \in [0, \infty).$$

If the weight distributions F_1, F_2, \dots, F_J are all typical (see Definition 2.1), then the mixture distribution \tilde{F} is also typical (see Section 2.4), and Theorem 2.1 immediately applies.

In contrast, if item types are revealed upon arrival, then the decision maker could use the type information to make better decisions. As we will see shortly, because the rewards are all equal, knowing the weight type of the arriving item makes no difference. The offline solution is still given by an algorithm that sorts items according to their realized weights (regardless of their types), and the optimal dynamic programming policy is a threshold policy that ignores weight types.

For the optimal offline solution, we can reinterpret this formulation so that items arrive according to a Bernoulli process with arrival probability $1 - p_0 = \sum_{j \in [J]} p_j$, have rewards equal to r and independent weights with distribution given by \tilde{F} . The optimal offline solution $R_n^*(c, 1 - p_0, r)$ is

then given by the sorting algorithm (2.1), so if

$$\epsilon_{k(1-p_0)}(x) = \sup \left\{ \epsilon \in [0, \infty) : \int_0^\epsilon w d\tilde{F}(w) \leq \frac{x}{k(1-p_0)} \right\}, \quad (2.37)$$

then Proposition 2.1 gives us that

$$\mathbb{E}[R_n^*(c, 1-p_0, r)] \leq n(1-p_0)r\tilde{F}(\epsilon_{n(1-p_0)}(c)), \quad (2.38)$$

and the prophet upper bound for weight distribution with multiple types follows.

To establish the independence on weight types of the optimal online solution when the rewards are all equal, we now examine the associated Bellman equation. We suppose that, with k periods to the end of the horizon, the remaining capacity is $x \in [0, c]$, the arriving item has weight type $j \in \{0, 1, \dots, J\}$ (with $j = 0$ denoting a no arrival or, equivalently, an arrival with arbitrarily large weight), and we let $V_k(x, j)$ be the optimal expected reward to-go given the current state. The optimality principle of dynamic programming then tells us that the value function $V_k(x, j)$ satisfies the Bellman recursion

$$V_k(x, j) = (1 - F_j(x)) \sum_{\iota=0}^J p_\iota V_{k-1}(x, \iota) + \int_0^x \max \left\{ r + \sum_{\iota=0}^J p_\iota V_{k-1}(x - w, \iota), \sum_{\iota=0}^J p_\iota V_{k-1}(x, \iota) \right\} dF_j(w), \quad (2.39)$$

with the initial condition $V_0(x, j) = 0$ for all $x \in [0, c]$ and all $j \in \{0, 1, \dots, J\}$. Here, the first summand holds because with probability $1 - F_j(x)$ the arriving type- j item has weight that exceeds the current knapsack capacity and the decision maker must reject it. Thus, her expected reward to-go over the remaining $k - 1$ periods is just given by the average over types of the value functions $V_{k-1}(x, \iota)$ for $\iota \in \{0, 1, \dots, J\}$. Instead, with probability $F_j(x)$ the arriving type- j item can be selected and the decision maker chooses the action that yields the largest expected reward to-go. If the item has weight w then its selection yields $r + \sum_{\iota=0}^J p_\iota V_{k-1}(x - w, \iota)$, while its rejection gives $\sum_{\iota=0}^J p_\iota V_{k-1}(x, \iota)$. By integrating this against $F_j(\cdot)$ for $w \in [0, x]$, we obtain the second summand of (2.39). The value functions $V_k(x, j)$ are monotone increasing in x for each j and k , and one has that

$$H_k^*(x, j) = \sup \left\{ w \in [0, x] : r + \sum_{\iota=0}^J p_\iota V_{k-1}(x - w, \iota) \geq \sum_{\iota=0}^J p_\iota V_{k-1}(x, \iota) \right\},$$

is the optimal threshold that identifies the largest type- j weight that can be selected when the current capacity is x and there are k periods remaining. Interestingly, one immediately has that

$H_k^*(x, j) = H_k^*(x, \iota)$ for all $j, \iota \in \{0, 1, \dots, J\}$ since all items have the same reward r and the expected rewards to-go of both actions are type independent. Because the optimal threshold policy ignores types, we can construct a heuristic that has the same property and use our earlier analysis to assess its performance. We recall the quantity $\epsilon_{k(1-p_0)}(x)$ in (2.37) and consider the type-independent threshold

$$\hat{H}_k(x, j) = \min\{x, \epsilon_{k(1-p_0)}(x)\} \quad \text{for all } j \in \{0, 1, \dots, J\} \text{ and } x \in [0, c].$$

If $\hat{\pi}$ is the heuristic that uses the thresholds $\hat{H}_n, \hat{H}_{n-1}, \dots, \hat{H}_1$, and $R_n^{\hat{\pi}}(c, 1 - p_0, r)$ is the total reward that $\hat{\pi}$ collects, then Proposition 2.3 tells us that there is a constant $1 < M < \infty$ depending only on \tilde{F} , the arrival probability $1 - p_0$, and the reward r such that

$$n(1 - p_0)r\tilde{F}(\epsilon_{n(1-p_0)}(c)) - M \log n \leq \mathbb{E} \left[R_n^{\hat{\pi}}(c, 1 - p_0, r) \right]. \quad (2.40)$$

If we combine the two bounds (2.38) and (2.40), we then have the corollary below.

Corollary 2.1 (Regret bound for weight distributions with multiple types). *Consider a knapsack problem with capacity $0 \leq c < \infty$ and with items that arrive over $1 \leq n < \infty$ periods according to a multinomial process with parameters $\mathbf{p} \equiv (p_0, p_1, \dots, p_J)$ such that $1 - p_0 = \sum_{j \in [J]} p_j$, and where p_0 is the probability of no arrival. If the items have rewards all equal to r and type-dependent weights with continuous distributions F_1, F_2, \dots, F_J and mixture (conditional on an arrival occurring) given by*

$$\tilde{F}(w) = \frac{1}{1 - p_0} \sum_{j \in [J]} p_j F_j(w) \quad \text{for all } w \in [0, \infty),$$

then

$$\mathbb{E} [R_n^*(c, 1 - p_0, r)] \leq n(1 - p_0)r\tilde{F}(\epsilon_{n(1-p_0)}(c)).$$

Furthermore, there is a feasible online policy $\hat{\pi}$ such that if the weights are independent and their distributions F_1, \dots, F_J belong to the typical class then there is a constant M depending only on \tilde{F} , p_0 , and r for which

$$n(1 - p_0)r\tilde{F}(\epsilon_{n(1-p_0)}(c)) - M \log n \leq \mathbb{E} \left[R_n^{\hat{\pi}}(c, 1 - p_0, r) \right].$$

In turn, if the weights are independent and F_1, \dots, F_J all belong to the typical class, then we have the regret bound

$$\mathbb{E} [R_n^*(c, 1 - p_0, r)] - \mathbb{E} \left[R_n^{\hat{\pi}}(c, 1 - p_0, r) \right] \leq M \log n.$$

We note here that the key assumption that makes our analysis carry over to weight distributions with multiple types is that the rewards are all equal. If one were to allow for type-dependent rewards, then the optimal offline solution would *not* be given by the offline-sort algorithm (2.1) and the optimal online solution would *not* be given by type-independent thresholds. While one would still have a Bellman recursion analogous to (2.39), it is unclear how type-dependent rewards would affect our regret estimates, and we leave this interesting open problem for future research.

2.8 Concluding Remarks

In this chapter we study the dynamic and stochastic knapsack problem with equal rewards and independent random weights with common continuous distribution F . We prove that—under some mild regularity conditions on the weight distribution—the regret is, at most, logarithmic in n . In particular, we show that this regret bound is attained by a reoptimized heuristic that can be expressed explicitly and that provides a key analytical connection with the offline solution.

Two questions stem naturally from our analysis. The first one entails the difference in performance between the reoptimized heuristic and the optimal online policy. Based on our numerical experiments, we conjecture that

$$\max_{\pi \in \Pi(n, c, p)} \mathbb{E}[R_n^\pi(c, r, p)] = \mathbb{E}[R_n^{\hat{\pi}}(c, r, p)] + O(1) \quad (2.41)$$

for all $n \geq 1$ and for a large class of weight distributions. However, it is well-known that the optimal policy often lacks of desirable structural properties, so proving (2.41) is unlikely to be easy. The second question has to do with the performance of the offline-sort algorithm. Here, the numerical evidence suggests that

$$\mathbb{E}[R_n^*(c, r, p)] = nprF(\epsilon_{np}(c)) + O(1)$$

for all $n \geq 1$ and most continuous weight distributions F .

Resolving the two conjectures above would imply that the regret cannot be $o(\log n)$ as $n \rightarrow \infty$ for most continuous weight distributions, and that $O(\log n)$ as $n \rightarrow \infty$ correctly quantifies the informational advantage that the prophet has over the sequential decision maker. This is in contrast with some other dynamic and stochastic knapsack problems in which the sequential decision maker does essentially as well as the prophet (see Section 2.1). It also suggests that when items have random weights, then the design of near-optimal heuristics requires more care than usual.

Chapter 3

An Adaptive $O(\log n)$ -Optimal Policy for the Sequential Selection of a Monotone Subsequence from a Random Sample

In the problem of *sequential* (online) selection of a *monotone increasing* subsequence, a decision maker observes sequentially a sequence of independent non-negative random variables X_1, X_2, \dots with common continuous distribution F and seeks to construct a monotone subsequence

$$X_{\tau_1} \leq X_{\tau_2} \leq \dots \leq X_{\tau_j} \quad (3.1)$$

where the indices $1 \leq \tau_1 < \tau_2 < \dots < \tau_j$ are stopping times with respect to the σ -fields $\mathcal{F}_i = \sigma\{X_1, X_2, \dots, X_i\}$, $1 \leq i < \infty$, and the trivial σ -field \mathcal{F}_0 . Since the indices are required to be possible values of stopping times, all selection/rejection decisions are terminal. That is, if the decision maker chooses not to select the value X_i at time i , then that value is lost forever. Similarly, if X_i is selected at time i , then that selection cannot be changed in the future. In general, the stopping times can be chosen to optimize different objective functions, and two main approaches have been considered in the literature. In the first, the decision maker seeks to maximize the expected number of selected elements when n are sequentially revealed (Samuels and Steele, 1981). In contrast, in the second approach the decision maker's objective is to minimize the expected time it takes to construct a monotone subsequence with n elements (Arlotto et al., 2016). Here, we confine our attention to the first — more classical — approach.

We then call a sequence of stopping times $1 \leq \tau_1 < \tau_2 < \dots < \tau_j \leq n$ such that (3.1) holds a *feasible policy*, and we denote the set of all such policies by $\Pi(n)$. For any $\pi \in \Pi(n)$, we then let $L_n(\pi)$ be the random variable that counts the number of selections made by policy π for the sample $\{X_1, X_2, \dots, X_n\}$. That is,

$$L_n(\pi) = \max\{j : X_{\tau_1} \leq X_{\tau_2} \leq \dots \leq X_{\tau_j} \text{ where } 1 \leq \tau_1 < \tau_2 < \dots < \tau_j \leq n\}.$$

Samuels and Steele (1981) first studied this selection problem and found that for each $n \geq 1$ there

This chapter is written under the supervision of Prof. Alessandro Arlotto and Prof. Yehua Wei. The results presented here are also in a joint paper Arlotto, Wei, and Xie (2018), published in *Random Structures & Algorithms*.

is a unique policy $\pi^* \in \Pi(n)$ such that

$$\mathbb{E}[L_n(\pi^*)] = \sup_{\pi \in \Pi(n)} \mathbb{E}[L_n(\pi)], \quad (3.2)$$

and for such optimal policies one has the asymptotic estimate

$$\mathbb{E}[L_n(\pi^*)] \sim (2n)^{1/2} \quad \text{as } n \rightarrow \infty. \quad (3.3)$$

Over the last few decades the understanding of policy π^* has substantially evolved. For instance, by formulating (3.2) as a finite-horizon Markov decision problem one sees that the optimal policy π^* is characterized by time and state dependent acceptance intervals. Furthermore, the asymptotic estimate (3.3) was refined with the much tighter bounds

$$(2n)^{1/2} - O(\log n) \leq \mathbb{E}[L_n(\pi^*)] \leq (2n)^{1/2} \quad \text{as } n \rightarrow \infty. \quad (3.4)$$

The upper bound in (3.4) was first discovered by Bruss and Robertson (1991) while studying the maximal number of elements of a random sample whose sum is less than a specified value. The analysis was instigated by the work of Coffman et al. (1987), and it now represents one of the early steps into the domain of resource-dependent branching processes (see, e.g., Bruss and Duerinckx, 2015). The result of Bruss and Robertson (1991) is actually quite rich; recent extensions and applications are discussed in Steele (2016). The upper bound in (3.4) also appeared in Gneden (1999) who considered the sequential selection of a monotone increasing subsequence from a random sample with random size.

The $O(\log n)$ lower bound in (3.4) is much more recent. It first appeared in the work of Bruss and Delbaen (2001) who studied the mean-optimal sequential selection of a monotone increasing subsequence when the observations X_1, X_2, \dots are revealed at the arrival epochs of a unit-rate Poisson process on $[0, n]$. While the Bruss and Delbaen (2001) result provides compelling evidence that a similar bound should also hold for the discrete-time formulation we consider here — i.e. the formulation in which the observations are revealed at the times $1, 2, \dots, n$ — the sequential nature of the two selection processes makes the result of Bruss and Delbaen (2001) not immediately applicable. The connection between the continuous-time formulation of Bruss and Delbaen (2001) and the discrete-time optimization (3.2) was then argued by Arlotto et al. (2015) who used the concavity of the map $n \mapsto \mathbb{E}[L_n(\pi^*)]$ and the $O(\log n)$ -bound of Bruss and Delbaen (2001) to ultimately confirm the lower bound in (3.4).

After a careful analysis, Bruss and Delbaen (2004) proved that the mean-optimal number of monotone increasing selections with Poisson-many observations is asymptotically normal after centering around $(2n)^{1/2}$ and scaling by $3^{-1/2}(2n)^{1/4}$. Arlotto et al. (2015) showed that the same asymptotic limit also holds for the discrete-time problem with n observations so, in summary, we now know that

$$\frac{3^{1/2}\{L_n(\pi^*) - (2n)^{1/2}\}}{(2n)^{1/4}} \implies N(0, 1), \quad \text{as } n \rightarrow \infty. \quad (3.5)$$

However, the analyses of Bruss and Delbaen (2001, 2004) and Arlotto et al. (2015) do not address whether there is a simple *adaptive* online policy — i.e., a policy that depends on the value of the last selection and on the number of observations that are yet to be seen — that is $O(\log n)$ optimal. The works of Rhee and Talagrand (1991) and Arlotto and Steele (2011) tell us that the best non-adaptive policy is $O(n^{1/4})$ optimal, but this optimality gap is too crude. For instance, the expected number of monotone increasing selections made by the best non-adaptive policy cannot even be used to center the random variable $L_n(\pi^*)$ around $(2n)^{1/2}$ in the weak law (3.5).

In this chapter, we construct a *simple adaptive online policy* $\hat{\pi}$ that is $O(\log n)$ optimal. The policy is characterized by a sequence of functions $\hat{\eta}_n, \hat{\eta}_{n-1}, \dots, \hat{\eta}_1$ such that if the value of the last selection up to and including time i is s and if $k = n - i$ observations remain to be seen then the value X_{i+1} is selected if and only if X_{i+1} falls in the *acceptance interval* $[s, \hat{\eta}_{n-i}(s)]$. In terms of the stopping times, the policy $\hat{\pi}$ corresponds to setting $\hat{\tau}_0 = 0$, $X_{\hat{\tau}_0} = 0$, and then defining the stopping times $\hat{\tau}_1 < \hat{\tau}_2 < \dots < \hat{\tau}_j$ recursively as

$$\hat{\tau}_j = \min\{\hat{\tau}_{j-1} < i \leq n : X_i \in [X_{\hat{\tau}_{j-1}}, \hat{\eta}_{n-i+1}(X_{\hat{\tau}_{j-1}})]\} \quad \text{for } 1 \leq j \leq n,$$

with the convention that if the set of indices on the right-hand side is empty, then $\hat{\tau}_j = \infty$. The random variable $L_n(\hat{\pi})$ then denotes the number of monotone increasing selections made by policy $\hat{\pi}$, and the expected value of $L_n(\hat{\pi})$ satisfies the two bounds given in the next theorem.

Theorem 3.1 ($O(\log n)$ -Optimal Policy). *For each $n \geq 1$, there is a simple adaptive online policy $\hat{\pi}$ such that*

$$(2n)^{1/2} - 2\{\log(n) + 1\} \leq \mathbb{E}[L_n(\hat{\pi})] \leq \mathbb{E}[L_n(\pi^*)] \leq (2n)^{1/2}. \quad (3.6)$$

We discuss the structure of policy $\hat{\pi}$ and of the functions $\{\hat{\eta}_k : 1 \leq k < \infty\}$ in Section 3.1, and then we turn to the proof of Theorem 3.1. The upper bound in (3.6) immediately follows from

(3.4) and the optimality of policy π^* , but our analysis requires a generalization of the upper bound (3.4) which we study in Section 3.2. Then in Section 3.3, we establish the equivalence between the sequential monotone subsequence selection problem and a special instance of the dynamic and stochastic knapsack problem with equal rewards of Chapter 2. With this equivalence, we can directly apply our previous results to prove the lower bound (3.6). Finally, in Section 3.4 we make concluding remarks and underscore some open problems.

3.1 Policy $\hat{\pi}$ and Its Value Function

For any feasible policy π and any continuous distribution F , we see that the number of selections made by π for the sample $\{X_1, X_2, \dots, X_n\}$ is unchanged if we replace each X_i by its monotone transformation $F^{-1}(X_i)$. Thus, we can assume without loss of generality that the X_i 's are uniformly distributed on $[0, 1]$. Next, we let

$$\hat{\eta}_k(s) = \min \left\{ s + [2k^{-1}(1-s)]^{1/2}, 1 \right\} \quad \text{for all } s \in [0, 1] \text{ and all } k \geq 1, \quad (3.7)$$

and we use the sequence of functions $\{\hat{\eta}_k : 1 \leq k < \infty\}$ to construct appropriate acceptance intervals. Specifically, if s denotes the value of the last selection when k observations are yet to be seen and x is the k th-to-last presented value, then x is selected as element of the subsequence that is under construction if and only if $x \in [s, \hat{\eta}_k(s)]$.

We now define the *critical value* $s_k = \max\{1 - 2k^{-1}, 0\}$ and we note that for $s \in [0, s_k]$ the decision maker is *conservative* and selects the k th-to-last value x if and only if it is within $\{2k^{-1}(1-s)\}^{1/2}$ of the most recent selection s . On the other hand, if $s \in [s_k, 1]$ the decision maker is *greedy* and accepts any k th-to-last value x that is larger than the most recent selection s .

If k denotes the number of observations that are yet to be seen and s is the value of the most recent selection, then we let $\hat{\nu}_k(s)$ denote the expected number of monotone increasing selections made by the acceptance interval policy characterized by the functions $\hat{\eta}_k, \hat{\eta}_{k-1}, \dots, \hat{\eta}_1$. The functions $\{\hat{\nu}_k : 1 \leq k \leq n\}$ are the value functions associated with policy $\hat{\pi}$ and they can be obtained recursively. Specifically, if $\hat{\nu}_0(s) = 0$ for all $s \in [0, 1]$, then for $k \geq 1$ we have the recursion

$$\hat{\nu}_k(s) = \{1 - \hat{\eta}_k(s) + s\}\hat{\nu}_{k-1}(s) + \int_s^{\hat{\eta}_k(s)} \{1 + \hat{\nu}_{k-1}(x)\} dx. \quad (3.8)$$

To see why this recursion holds, we condition on the k th-to-last uniform random value. With

probability $1 - \hat{\eta}_k(s) + s$ the newly presented value x does not fall in the acceptance interval $[s, \hat{\eta}_k(s)]$. In this case no selection is made and we are left with $k - 1$ remaining observations and with the value of the most recent selection s unchanged. This amounts to an expected number of remaining selections equal to $\hat{\nu}_{k-1}(s)$ and it justifies the first summand of our recursion (3.8). On the other hand, with probability $\hat{\eta}_k(s) - s$ the newly presented value x falls in the acceptance interval, and we obtain a reward of one for selecting x plus the expected number of remaining selections over the next $k - 1$ observations when the value of the most recent selection changes to x . Integrating this over all $x \in [s, \hat{\eta}_k(s)]$ gives us the second summand of the recursive equation (3.8). The value functions $\{\hat{\nu}_k : 1 \leq k < \infty\}$ are all continuous on $[0, 1]$ and their behavior is well summarized by the following theorem.

Theorem 3.2 ($\hat{\nu}_k$ Bounds). *For all $k \geq 1$ and all $s \in [0, 1]$ one has that*

$$\{2k(1 - s)\}^{1/2} - 2\{\log(k) + 1\} \leq \hat{\nu}_k(s) \leq \{2k(1 - s)\}^{1/2}. \quad (3.9)$$

Since policy $\hat{\pi}$ is characterized by the thresholds $\hat{\eta}_n, \hat{\eta}_{n-1}, \dots, \hat{\eta}_1$ and by the initial state $s = 0$, one then has the equivalence

$$\mathbb{E}[L_n(\hat{\pi})] = \hat{\nu}_n(0) \quad \text{for all } n \geq 1.$$

Hence, Theorem 3.1 is an immediate corollary of Theorem 3.2, but the upper bound (3.9) is a refinement of (3.4) to an arbitrary initial state. In fact the same upper bound holds for all feasible policies based on acceptance intervals, including the optimal one.

The estimates in Theorem 3.2 also allow us to provide some intuition for our choice of the threshold functions $\{\hat{\eta}_k : 1 \leq k < \infty\}$ in (3.7). For every $k \geq 1$ and $s \in [0, 1]$ the threshold functions aim to balance the expected reward to-go $\hat{\nu}_{k-1}(s)$ that one obtains when skipping the k th-to-last observation x , and the reward $1 + \hat{\nu}_{k-1}(x)$ that one earns when selecting the k th-to-last value x . Since $\hat{\nu}_{k-1}(s) \approx \{2(k-1)(1-s)\}^{1/2}$, one can solve the equation

$$\{2(k-1)(1-s)\}^{1/2} = 1 + \{2(k-1)(1-\hat{x})\}^{1/2} \quad (3.10)$$

to find the largest value of x that makes selecting the current value worthwhile. Equation (3.10) then tells us that

$$\hat{x} = s + [2(k-1)^{-1}(1-s)]^{1/2} - [2(k-1)]^{-1} \approx s + [2k^{-1}(1-s)]^{1/2} - [2k]^{-1},$$

so our choice (3.7) accounts for the first two terms of the approximation of the solution of equation (3.10). The truncation in (3.7) then ensures that all the thresholds $\{\hat{\eta}_k : 1 \leq k < \infty\}$ are feasible. At the end of the next section we use an optimization argument to provide further intuition for our choice of the threshold functions.

3.2 A Refined Prophet Upper Bound

In this section, we prove the upper bound (3.9) by showing that it holds for all policies that are based on acceptance intervals. The adaptive policy $\hat{\pi}$ and the unique optimal policy π^* both have this property. The argument we provide draws substantially from the earlier analyses of Gnedin (1999) and Bruss and Delbaen (2001), but it takes advantage of the flexibility that comes from allowing for an arbitrary initial state. Specifically, here we assume that the first subsequence element can be selected only if it is larger than an arbitrary value $s \in [0, 1]$. In contrast, in the classical formulation one always takes the initial state $s = 0$.

For any $k \geq 1$, an arbitrary acceptance interval policy π_k is given by a sequence $\eta_k, \eta_{k-1}, \dots, \eta_1$ of functions such that

$$s \leq \eta_j(s) \leq 1 \quad \text{for all } 1 \leq j \leq k \text{ and all } s \in [0, 1].$$

If s denotes the initial state or the value of the last observation selected prior to being presented with the j -to-last value x , then x is selected if and only if $x \in [s, \eta_j(s)]$. Next, for any $s \in [0, 1]$ we set $M_0 = s$ and we let M_i denote the maximum between M_0 and the largest of the elements of the subsequence that have been selected up to and including time i . The number of selections from $\{X_1, X_2, \dots, X_k\}$ when $M_0 = s$ is then given by the random variable

$$L_k(\pi_k, s) = \sum_{i=1}^k \mathbb{1} \{X_i \in [M_{i-1}, \eta_{k-i+1}(M_{i-1})]\}.$$

If we now take expectations on both sides and use the Cauchy-Schwarz inequality to estimate the sum of the products $1 \cdot \mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}]$ for $1 \leq i \leq k$, we obtain

$$\begin{aligned} \mathbb{E}[L_k(\pi_k, s)] &= \sum_{i=1}^k 1 \cdot \mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}] \\ &\leq k^{1/2} \left\{ \sum_{i=1}^k \mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}]^2 \right\}^{1/2}. \end{aligned} \tag{3.11}$$

The definition of M_i as the maximum between M_0 and the largest of the elements of the subsequence that have been selected up to and including time i tells us that

$$M_i = \begin{cases} M_{i-1} & \text{if } X_i \notin [M_{i-1}, \eta_{k-i+1}(M_{i-1})] \\ X_i & \text{if } X_i \in [M_{i-1}, \eta_{k-i+1}(M_{i-1})] \end{cases}$$

and, because X_i is uniformly distributed on the unit interval, we have the identity

$$\mathbb{E}[M_i - M_{i-1} \mid \mathcal{F}_{i-1}] = \int_{M_{i-1}}^{\eta_{k-i+1}(M_{i-1})} (x - M_{i-1}) dx = \frac{1}{2} (\eta_{k-i+1}(M_{i-1}) - M_{i-1})^2.$$

Taking the total expectation then gives

$$\mathbb{E}[(\eta_{k-i+1}(M_{i-1}) - M_{i-1})^2] = 2\{\mathbb{E}[M_i] - \mathbb{E}[M_{i-1}]\},$$

so a second application of the Cauchy-Schwarz inequality implies the upper bound

$$\mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}]^2 \leq 2\{\mathbb{E}[M_i] - \mathbb{E}[M_{i-1}]\}.$$

If we now sum over $1 \leq i \leq k$ and recall that $\mathbb{E}[M_k] \leq 1$ and $M_0 = s$, we obtain from telescoping that

$$\sum_{i=1}^k \mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}]^2 \leq 2\{\mathbb{E}[M_k] - \mathbb{E}[M_0]\} \leq 2(1-s). \quad (3.12)$$

This last inequality and the bound in (3.11) finally give us that

$$\mathbb{E}[L_k(\pi_k, s)] \leq \{2k(1-s)\}^{1/2} \quad \text{for all } k \geq 1 \text{ and all } s \in [0, 1]$$

and complete the proof of the upper bound (3.9).

The same upper bound can also be obtained by formulating a simple optimization problem that provides further insight into our choice of the adaptive thresholds $\{\hat{\eta}_k : 1 \leq k < \infty\}$. We consider the optimization problem

$$\begin{aligned} w^* &= \max_{d_1, \dots, d_k} \sum_{i=1}^k d_i \\ \text{s.t. } & \sum_{i=1}^k d_i^2 \leq 2(1-s), \end{aligned} \quad (3.13)$$

and we obtain from inequality (3.12) that $d_i = \mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}]$ is a feasible solution. This feasible solution has objective

$$\sum_{i=1}^k \mathbb{E}[\eta_{k-i+1}(M_{i-1}) - M_{i-1}] = \mathbb{E}[L_k(\pi_k, s)],$$

so we have that the optimal objective value of (3.13) is an upper bound for $\mathbb{E}[L_k(\pi_k, s)]$, and that (3.13) is a relaxation of the sequential monotone subsequence problem (3.2).

The optimal value of (3.13) can be easily estimated. The problem has a linear objective function and convex constraints; by the Karush–Kuhn–Tucker (KKT) conditions, its optimal solution $(d_1^*, d_2^*, \dots, d_k^*)$ is given by

$$d_i^* = [2k^{-1}(1-s)]^{1/2} \quad \text{for all } 1 \leq i \leq k,$$

and its optimal value is

$$w^* = \{2k(1-s)\}^{1/2} \quad \text{for all } k \geq 1 \text{ and } s \in [0, 1].$$

It then follows that the adaptive thresholds $\{\hat{\eta}_k : 1 \leq k < \infty\}$ defined by (3.7) are *reoptimized* thresholds that use the optimal solution of the relaxation (3.13) for any given $k \geq 1$ and $s \in [0, 1]$. Of course, a small nuisance arises, and one needs to make sure that the reoptimized thresholds are indeed feasible. Since there are values $s \in [0, 1]$ such that $s + [2k^{-1}(1-s)]^{1/2} > 1$, one needs the truncation introduced in (3.7) to obtain thresholds feasible. In the next section, we connect the sequential monotone subsequence selection problem to the dynamic and stochastic knapsack problem of Chapter 2, and apply the established results to prove the lower bound in Theorem 3.2.

3.3 Equivalence with the Dynamic and Stochastic Knapsack Problem with Equal Rewards

As we mentioned in Chapter 2, the problem of the sequential selection of a monotone subsequence from a random sample is equivalent to a special instance of the dynamic and stochastic knapsack problem with equal rewards. More specifically, the equivalent special instance is the one in which the items arrive deterministically in every period, have unitary rewards and independent random weights drawn from the uniform distribution on $[0, 1]$, and the knapsack has unit capacity. Such equivalence can be easily seen with one change of variable once one writes down the Bellman

equations for the two problems. Specifically, for the problem of sequential increasing subsequence selection, the Bellman equation is given by the recursion

$$\nu_k^*(s) = s\nu_{k-1}^*(s) + \int_s^1 \max\{1 + \nu_{k-1}^*(x), \nu_{k-1}^*(s)\} dx, \quad \text{for } k \geq 1, \quad (3.14)$$

and by the boundary condition $\nu_0^*(s) \equiv 0$ for all $s \in [0, 1]$. The value function $\nu_k^*(s)$ represents the expected number of monotone increasing selections made by the optimal policy when there are k values yet to be seen and the last selected subsequence element is s . The optimality of the Bellman equation (3.14) can be verified by conditioning on the k th-to-last uniform random value x . On the one hand, with probability s this newly realized value is smaller than the last selected value s . In this case, no selection can be made and one is left with $k - 1$ random values yet to observe and the last selected value s is unchanged. This yields the first term on the right-hand side of (3.14). On the other hand, if the uniform random value $x \geq s$, then one chooses to select or to reject this value. A selection of x yields one element in the selected subsequence and an expected number of future selections $\nu_{k-1}^*(x)$ since there are $k - 1$ random values and the last selected value now is x . On the contrary, a rejection yields no selected element and an expected number of future selections $\nu_{k-1}^*(s)$ since the last selected value s is unchanged. By integrating the maximum of the two yields over the interval $[s, 1]$, we recover the second term on the right-hand side of the optimality equation (3.14).

If we now consider a dynamic and stochastic knapsack problem with initial capacity $c = 1$, items that arrive in each period with probability $p = 1$, rewards $r = 1$, and random weights that are drawn uniformly on the unit interval, then we obtain from (2.36) that the Bellman equation is given by the recursion

$$v_k^*(s) = (1 - s)v_{k-1}^*(s) + \int_0^s \max\{1 + v_{k-1}^*(s - w), v_{k-1}^*(s)\} dw, \quad \text{for } k \geq 1, \quad (3.15)$$

and by the boundary condition $v_0^*(s) \equiv 0$ for all $s \in [0, 1]$.

The equivalence between the above two problems is summarized in the following proposition.

Proposition 3.1. *For all $k \geq 1$ and all $s \in [0, 1]$, one has that*

$$\nu_k^*(s) = v_k^*(1 - s).$$

Proof. The two recursions have the same boundary condition $\nu_0^*(s) = v_0^*(1 - s) = 0$ for all $s \in [0, 1]$,

so we have the induction case. Now, we take as induction hypothesis that

$$\nu_{k-1}^*(s) = v_{k-1}^*(1-s) \quad \text{for all } s \in [0, 1],$$

and we prove that $\nu_k^*(s) = v_k^*(1-s)$ for all $s \in [0, 1]$. By replacing s with $1-s$ in (3.15) we obtain that

$$v_k^*(1-s) = (1-s)v_{k-1}^*(1-s) + \int_0^{1-s} \max\{1 + v_{k-1}^*(1-s-w), v_{k-1}^*(1-s)\} dw.$$

On the right-hand side of this last equation, we can now apply the induction hypothesis to replace $v_{k-1}^*(1-s)$ with $\nu_{k-1}^*(s)$ and to replace $v_{k-1}^*(1-s-w)$ with $\nu_{k-1}^*(s+w)$. We then obtain that

$$v_k^*(1-s) = (1-s)\nu_{k-1}^*(s) + \int_0^{1-s} \max\{1 + \nu_{k-1}^*(s+w), \nu_{k-1}^*(s)\} dw.$$

If we now change variable and let $w = x - s$ to this last right-hand side, we obtain that

$$v_k^*(1-s) = (1-s)\nu_{k-1}^*(s) + \int_s^1 \max\{1 + \nu_{k-1}^*(x), \nu_{k-1}^*(s)\} dx.$$

By recalling the Bellman equation (3.14) for the sequential monotone subsequence selection problem, we finally obtain that $\nu_k^*(s) = v_k^*(1-s)$ for all $s \in [0, 1]$, just as needed. \square

In addition to the equivalence between the optimal policies of the two problems, one also has the equivalence between the adaptive policy $\hat{\pi}$ of Section 2.3 and the adaptive policy $\hat{\pi}$ of Section 3.1. To see this equivalence, let us recall that for the special instance of the dynamic and stochastic knapsack problem, we know from (2.2) and (2.18) that the threshold of policy $\hat{\pi}$ for the knapsack problem is given by

$$\hat{h}_k(s) = \min\{s, (2k^{-1}s)^{1/2}\} \quad \text{for all } s \in [0, 1] \text{ and all } k \geq 1. \quad (3.16)$$

On the other hand, the threshold of policy $\hat{\pi}$ for the sequential monotone subsequence selection problem in (3.7) is given by

$$\hat{\eta}_k(s) = \min\left\{s + [2k^{-1}(1-s)]^{1/2}, 1\right\} \quad \text{for all } s \in [0, 1] \text{ and all } k \geq 1, \quad (3.17)$$

and one can easily see that the two thresholds (3.16) and (3.17) satisfy the following relation

$$\hat{h}_k(1-s) = \hat{\eta}_k(s) - s \quad (3.18)$$

The identity above provides the essential link to connect the value function of the policy $\hat{\pi}$ for the dynamic and stochastic knapsack problem in (2.19) and given by

$$\hat{v}_k(s) = (1 - \hat{h}_k(s))\hat{v}_{k-1}(s) + \int_0^{\hat{h}_k(s)} (1 + \hat{v}_{k-1}(s - w)) dw, \quad (3.19)$$

with the value function of policy $\hat{\pi}$ for the sequential monotone subsequence selection problem, which we recall from (3.8), as

$$\hat{v}_k(s) = \{1 - \hat{\eta}_k(s) + s\}\hat{\nu}_{k-1}(s) + \int_s^{\hat{\eta}_k(s)} \{1 + \hat{\nu}_{k-1}(x)\} dx. \quad (3.20)$$

Just as in Proposition 3.1, we see that the two recursions (3.19) and (3.20) are equivalent in the following sense.

Proposition 3.2. *For all $k \geq 1$ and all $s \in [0, 1]$, one has that*

$$\hat{\nu}_k(s) = \hat{v}_k(1 - s).$$

Proof. The boundary condition $\hat{\nu}_0(s) = \hat{v}_0(1 - s) = 0$ for all $s \in [0, 1]$ gives us the induction case. Next, we prove that $\hat{\nu}_k(s) = \hat{v}_k(1 - s)$ with the induction hypothesis

$$\hat{\nu}_{k-1}(s) = \hat{v}_{k-1}(1 - s) \quad \text{for all } s \in [0, 1].$$

If we replace s with $1 - s$ in the recursion (3.19), we obtain that

$$\hat{v}_k(1 - s) = \{1 - \hat{h}_k(1 - s)\}\hat{v}_{k-1}(1 - s) + \int_0^{\hat{h}_k(1 - s)} \{1 + \hat{v}_{k-1}(1 - s - w)\} dw.$$

By invoking the induction hypothesis, we can replace $\hat{v}_{k-1}(1 - s)$ with $\hat{\nu}_{k-1}(s)$ and replace $\hat{v}_{k-1}(1 - s - w)$ with $\hat{\nu}_{k-1}(s + w)$ to obtain that

$$\hat{v}_k(1 - s) = \{1 - \hat{h}_k(1 - s)\}\hat{\nu}_{k-1}(s) + \int_0^{\hat{h}_k(1 - s)} \{1 + \hat{\nu}_{k-1}(s + w)\} dw.$$

Then if we perform a change of variable $w = x - s$ for the integral on this last right-hand side, we obtain that

$$\hat{v}_k(1 - s) = \{1 - \hat{h}_k(1 - s)\}\hat{\nu}_{k-1}(1 - s) + \int_s^{\hat{h}_k(1 - s) + s} \{1 + \hat{\nu}_{k-1}(x)\} dx.$$

Lastly, if we compare the two thresholds (3.16) and (3.17), and plug in the identity $\hat{h}_k(1-s) = \hat{\eta}_k(s) - s$ in (3.18), we obtain that

$$\hat{v}_k(1-s) = \{1 - \hat{\eta}_k(s) + s\} \hat{v}_{k-1}(1-s) + \int_s^{\hat{\eta}_k(s)} \{1 + \hat{v}_{k-1}(x)\} dx. \quad (3.21)$$

By (3.20) we know that this last right-hand side is equal to $\hat{v}_k(s)$, and hence $\hat{v}_k(s) = \hat{v}_k(1-s)$. As a result,

$$\hat{v}_k(s) = \hat{v}_k(1-s) \quad \text{for all } s \in [0, 1] \text{ and all } k \geq 1,$$

and the proof is complete. \square

With this last proposition, we can apply our previous analysis of the dynamic and stochastic knapsack problem in Chapter 2 to prove the lower bound in Theorem 3.2. In fact, Proposition 3.2 tells us that the value function $\hat{v}_k(s)$ of policy $\hat{\pi}$ for the sequential increasing subsequence selection problem is equal to the value function $\hat{v}_k(1-s)$ of the adaptive heuristic $\hat{\pi}$ of Section 2.3 for the dynamic and stochastic knapsack problem with equal rewards. As we discussed in Section 2.4, the uniform distribution on the interval $[0, 1]$ belongs to the typical class in Definition 2.1. Moreover, for this distribution, its associated constant in Lemma 2.2 is $M = 2$. Therefore, if we apply Theorem 2.1 to the value function $\hat{v}_k(1-s)$, we would obtain that for all $s \in [0, 1]$ and all $k \geq 1$,

$$\{2k(1-s)\}^{-1} - 2(\log(k) + 1) \leq \hat{v}_k(1-s).$$

Then we obtain from Proposition 3.2 that

$$\{2k(1-s)\}^{-1} - 2(\log(k) + 1) \leq \hat{v}_k(s) \quad \text{for all } s \in [0, 1] \text{ and all } k \geq 1,$$

completing the proof of the lower bound of (3.9) in Theorem 3.2.

We close this section with the comparison between the full-information version of the sequential monotone subsequence selection problem and that of the dynamic and stochastic knapsack problem. For a random sample $\{X_1, \dots, X_n\}$, the full-information version of the sequential increasing subsequence selection problem is to find the longest increasing subsequence with length

$$L_n \equiv \max \{j : X_{i_1} \leq X_{i_2} \leq \dots \leq X_{i_j} \text{ where } 1 \leq i_1 < i_2 < \dots < i_j \leq n\}. \quad (3.22)$$

This last solution is a sample-path upper bound for the sequential increasing subsequence selection problem. That is, for the optimal policy (3.2) of the sequential increasing subsequence selection

problem, one has that

$$L_n(\pi^*) \leq L_n \quad \text{almost surely.}$$

However, the upper bound $\mathbb{E}[L_n(\pi^*)] \leq (2n)^{1/2}$ in Theorem 3.1 does not apply to the full-information solution to the longest increasing subsequence selection problem. This does not come out very surprisingly because as we discussed in Section 3.2, such upper bound is intimately related to the fact that the optimal sequential policy is of a thresholding type. The solution (3.22) to the longest increasing subsequence selection problem, however, is not of a thresholding type due to its combinatorial nature. Moreover, for the sequential increasing subsequence selection problem, we know from Theorem 3.1 that

$$\mathbb{E}[L_n(\pi^*)] \sim (2n)^{1/2} \quad \text{as } n \rightarrow \infty.$$

While for the longest increasing subsequence selection problem, it has long been known that (Logan and Shepp, 1977; Vershik and Kerov, 1977)

$$\mathbb{E}[L_n] \sim 2(n)^{1/2} \quad \text{as } n \rightarrow \infty.$$

In other words, in the problems of monotone subsequence selection, a prophet with full information has a huge information advantage over the optimal sequential decision maker. This is in contrast with the dynamic and stochastic knapsack problem with equal rewards in which the online and offline solutions have the same first-order asymptotic performance. This is because both policies are given by a thresholding algorithm (Lemma 2.1). This feature is not present in the monotone subsequence selection problems in which only the online optimal algorithm is given by a sequence of threshold functions.

3.4 Connections and Observations

Policy $\hat{\pi}$ is a simple adaptive online policy that selects a monotone increasing subsequence. It is then easy to see how it could be generalized to the sequential selection of a unimodal or d -modal subsequence considered by Arlotto and Steele (2011). For instance, in the unimodal case with n observations, one could set the turning time $\bar{n} = \lfloor n/2 \rfloor$ and run policy $\hat{\pi}_{\bar{n}}$ to construct an increasing segment with the first \bar{n} observations, and a decreasing version of $\hat{\pi}_{n-\bar{n}}$ to obtain a decreasing segment over the next $n - \bar{n}$ observations. Theorem 3.1 then would immediately apply

to this construction, proving a $O(\log n)$ -optimality gap for the sequential selection of unimodal and d -modal subsequences.

It is also reasonable to expect that policy $\hat{\pi}$ could generalize to the sequential selection of coordinatewise increasing subsequences from a uniform random sample on the m -dimensional hypercube. Thus far, the best optimality-gap estimate comes from the work of Baryshnikov and Gnedin (2000) who use a non-adaptive policy to derive a $O(n^{1/(2m+2)})$ bound. The adaptive character of policy $\hat{\pi}$ could be fruitful one more time and help establish a $O(\log n)$ -optimality gap for this multidimensional problem.

In addition to the well-known bound $\mathbb{E}[L_n(\pi^*)] \leq (2n)^{1/2}$, Seksenbayev (2018) and Gnedin and Seksenbayev (2019) recently prove that the gap of this last upper bound has the asymptotics

$$g(n) = (2n)^{1/2} - \mathbb{E}[L_n(\pi^*)] \sim (12)^{-1} \log(n) \quad \text{as } n \rightarrow \infty.$$

This last asymptotics tells us that our bound of Theorem 3.1 on the optimality gap is order tight. We also note that in the closely related problem in which the sequence X_1, X_2, \dots, X_n is given by a uniform random permutation of the integers $\{1, 2, \dots, n\}$, Peng and Steele (2016) showed that the corresponding difference $g(n)$ is $O(\log n)$ as $n \rightarrow \infty$.

Several interesting questions remain open, however. Theorem 3.2 tells us that

$$\{2k(1-s)\}^{1/2} - \hat{v}_k(s) \leq 2\{\log(k) + 1\} \quad \text{for all } s \in [0, 1] \text{ and all } k \geq 1,$$

and it would be worthwhile to understand how the right-hand side changes with the initial state value s . When $s = 0$ we have from Theorem 3.1 that

$$\mathbb{E}[L_n(\pi^*)] \leq \mathbb{E}[L_n(\hat{\pi})] + 2\{\log(n) + 1\} \quad \text{for all } n \geq 1,$$

but the actual expected performance of policy $\hat{\pi}$ seems to be much tighter. Based on an extensive numerical analysis¹ (see also Figure 2.2), we conjecture that there is a constant $0 < c < \infty$ such that

$$\mathbb{E}[L_n(\pi^*)] \leq \mathbb{E}[L_n(\hat{\pi})] + c \quad \text{for all } n \geq 1.$$

Such conjecture is indeed a specialization of our early conjecture (2.41) of Section 2.8, and it would also imply that the functions $\hat{\eta}_n, \hat{\eta}_{n-1}, \dots, \hat{\eta}_1$ are remarkably close to their analogues that arise when implementing the optimal dynamic programming algorithm.

¹We estimated numerically the value functions that solve the recursion (3.8) and its optimal counterpart on a discretized state space with a grid size of 10^{-5} and with k ranging from 1 to 10^4 .

Chapter 4

Sequential Policies and the Distribution of Their Total Rewards in Dynamic and Stochastic Knapsack Problems

Originally introduced by Bellman (1957), stochastic dynamic programming is a technique for solving sequential decision problems in presence of uncertainty. Such problems can, in principle, be solved optimally by backward induction, but their solutions are rarely available in closed form, and they often suffer from the curse of dimensionality (Powell, 2011). As such, there is a pressing need for identifying approximate solutions with provable performance guarantees, and one common near-optimal criterion is that of asymptotic optimality (see, e.g., Coffman et al., 1987; Bruss and Robertson, 1991; Rhee and Talagrand, 1991; Gallego and van Ryzin, 1994, 1997; Talluri and van Ryzin, 1998). In this chapter, we study the reoptimized heuristic of Chapter 2 for the dynamic and stochastic knapsack problem with unitary rewards, and we show that the number of selections made by such policy has the same asymptotic variance and the same limiting distribution as that of the optimal dynamic programming policy. In contrast, we note that the number of selections made by another widely studied asymptotically optimal heuristic has larger asymptotic variance and a different limiting distribution. This contrast suggests that decision makers who are interested in assessing the higher-order effects of their near-optimal decisions should be careful with solely relying on asymptotic optimality.

To fix ideas and gain intuition, we now preview our main result by discussing a dynamic and stochastic knapsack problem with unitary rewards, and independent and identically distributed weights on $[0, 1]$. In such a problem, a decision maker (referred to as *she*) is given a knapsack with unit capacity, and she is presented with items arriving sequentially over n discrete time periods. As soon as an item arrives, its uniform weight $W_t, t \in [n] \equiv \{1, \dots, n\}$ is revealed, and the decision maker has to decide whether to include the item in the knapsack or not. All decisions are terminal; they cannot be changed at a later time. The decision maker's objective is to maximize the

This chapter is written under the supervision of Prof. Alessandro Arlotto. The results presented here are also in a joint research paper Arlotto and Xie (2020b).

expected number of selections under the capacity constraint, and this can be achieved by formulating the problem as a dynamic program, or by proposing near-optimal heuristics. For instance, if the remaining knapsack capacity at the beginning of period t is x , then Chapter 2 proposes the *adaptive* (time-dependent) heuristic $\hat{\pi}$ that selects the t -th arriving item if and only if its weight $W_t \leq \hat{h}_{n-t+1}(x) \equiv \min \left\{ x, \sqrt{\frac{2x}{n-t+1}} \right\}$. This decision rule then gives us a sequence of stopping times $0 \equiv \hat{\tau}_0 < \hat{\tau}_1 < \dots < \hat{\tau}_k$ such that

$$\hat{\tau}_j = \min \left\{ t : \hat{\tau}_{j-1} < t \leq n \text{ and } W_t \leq \hat{h}_{n-t+1} \left(1 - \sum_{i=1}^{j-1} W_{\hat{\tau}_i} \right) \right\} \quad \text{for } j \geq 1, \quad (4.1)$$

and a total number of item selections given by

$$\hat{N}_n = \max \left\{ j \in \{0, 1, \dots, n\} : \sum_{i=1}^j W_{\hat{\tau}_i} \leq 1 \right\}. \quad (4.2)$$

One can also choose alternative thresholds to make selection decisions. Bruss and Robertson (1991); Rhee and Talagrand (1991) suggest the *non-adaptive* policy $\tilde{\pi}$ that selects the t -th arriving item if and only if its weight $W_t \leq \tilde{h}_{n-t+1}(x) \equiv \min \left\{ x, \sqrt{\frac{2}{n}} \right\}$, where—just as before—the variable x denotes the remaining knapsack capacity. One then has a sequence of stopping times analogous to (4.1) with the associated number of item selections \tilde{N}_n .

Finally, one can also consider the thresholds h_n^*, \dots, h_1^* associated with the optimal dynamic programming policy π^* (Papastavrou et al., 1996) with its optimal number of item selections N_n^* . Here the three policies above play a key role because they all share the same first-order asymptotic performance.

In a sequence of papers (Samuels and Steele, 1981; Coffman et al., 1987; Rhee and Talagrand, 1991; Arlotto et al., 2018; Arlotto and Xie, 2020a), it has been shown that, all the three sequential policies have the same first-order performance with the closed-form expression

$$\mathbb{E} \left[\tilde{N}_n \right] \sim \mathbb{E} \left[\hat{N}_n \right] \sim \mathbb{E} \left[N_n^* \right] \sim \sqrt{2n} \quad \text{as } n \rightarrow \infty.$$

However, the second-order performances of the three policies are different. Specifically, Boshuizen and Kertz (1999) proves that the non-adaptive heuristic $\tilde{\pi}$ has the variance asymptotics

$$\text{Var} \left\{ \tilde{N}_n \right\} \sim \frac{2(\pi - 1)}{3\pi} \sqrt{2n} \quad \text{as } n \rightarrow \infty,$$

as well as the limiting behavior given as

$$\frac{\tilde{N}_n - \sqrt{2n}}{\sqrt[4]{2n}} \Rightarrow Z_1 + (Z_2 \wedge 0) \quad \text{as } n \rightarrow \infty,$$

where Z_1 and Z_2 are centered normal random variables with covariance matrix $\begin{pmatrix} 1 & -1 \\ -1 & 4/3 \end{pmatrix}$.

On the other hand, we will discover as corollary of the main result in this chapter that the other two policies have the same variance asymptotics

$$\text{Var} \left\{ \hat{N}_n \right\} \sim \text{Var} \left\{ N_n^* \right\} \sim \frac{1}{3} \sqrt{2n} \quad \text{as } n \rightarrow \infty,$$

which is smaller than that of the non-adaptive heuristic. In addition, with the same centering and scaling, our main result also implies the limit theorems

$$\frac{\hat{N}_n - \sqrt{2n}}{\sqrt[4]{2n}} \Rightarrow N \left(0, \frac{1}{3} \right) \quad \text{and} \quad \frac{N_n^* - \sqrt{2n}}{\sqrt[4]{2n}} \Rightarrow N \left(0, \frac{1}{3} \right) \quad \text{as } n \rightarrow \infty.$$

Several observations are in order (see also Table 4.1). First, all three policies are first-order equivalent since they have the same first-order asymptotics $\sqrt{2n}$. Second, the adaptive heuristic $\hat{\pi}$ shares the same variance asymptotics as the optimal dynamic programming policy π^* , while the non-adaptive heuristic $\tilde{\pi}$ has an asymptotic variance that is about 36% larger. As such, while policy $\tilde{\pi}$ provides the same asymptotic mean performance as policies $\hat{\pi}$ and π^* , it does so with much higher risk. Lastly, the adaptive heuristic $\hat{\pi}$ also shares the same limiting distribution as the optimal dynamic programming policy π^* , while the non-adaptive heuristic $\tilde{\pi}$ has a different limiting distribution. In other words, the adaptivity of $\hat{\pi}$ fully recovers the probabilistic behavior of the optimal dynamic programming policy π^* , going beyond asymptotic optimality.

As we will see in Section 4.1, such observations carry through from this special case to a class of weight distributions that satisfies certain regularity conditions. Specifically, we show that the adaptive heuristic $\hat{\pi}$ is equivalent to the optimal dynamic programming policy π^* . It provides the same asymptotic first-order performance, the same asymptotic second-order performance, and the same limiting distribution. In general, our results suggest that, when faced with sequential decision problems, decision makers should be well informed about the potential risks that come with choosing which heuristic to implement.

Table 4.1: Asymptotic performance comparison as $n \rightarrow \infty$ among three policies.

	Non-adaptive heuristic $\tilde{\pi}$	Adaptive heuristic $\hat{\pi}$	Optimal Dynamic programming policy π^*
$\mathbb{E}[N_n]$	$\sqrt{2n}$	$\sqrt{2n}$	$\sqrt{2n}$
$\text{Var}\{N_n\}$	$\frac{2(\pi-1)}{3\pi}\sqrt{2n}$	$\frac{1}{3}\sqrt{2n}$	$\frac{1}{3}\sqrt{2n}$
$\frac{N_n - \sqrt{2n}}{\sqrt[4]{2n}}$	$Z_1 + (Z_2 \wedge 0)$	$N\left(0, \frac{1}{3}\right)$	$N\left(0, \frac{1}{3}\right)$

Organization of the Chapter

In Section 4.1, we revisit the adaptive heuristic $\hat{\pi}$ of Chapter 2 and state our main result. In Section 4.2 we establish the key Lemma 4.1 to control the derivative $\hat{v}'_k(x)$ of the value function associated with the adaptive heuristic $\hat{\pi}$. This key lemma is then used in Section 4.3 to obtain the variance asymptotics of the number of selections made by the adaptive heuristic $\hat{\pi}$. With this last variance asymptotics and related estimations, we close our proof of the main result in Section 4.4. Lastly, we conclude this chapter with some closing remarks in Section 4.5.

4.1 Dynamic and Stochastic Knapsack Problem with Unitary Rewards: A Distributional Perspective

As we discussed in the literature review of Section 2.1, for the dynamic and stochastic knapsack problem and related formulations, much attention has been put to the analyses of policies in terms of their mean performances and worst-case performances. However, the distributional perspective of the problem has been much less touched in the literature. A notable exception is the case in which items arrive deterministically in each period with unit rewards, have uniform weights that are drawn on the interval $[0, 1]$, and the knapsack has unit capacity. For this case, the results of Arlotto et al. (2015) imply that the number of selections made by the optimal sequential policy, $N_n^*(1)$, has the variance asymptotics

$$\text{Var}\{N_n^*(1)\} \sim \frac{1}{3}(2n)^{1/2} \quad \text{as } n \rightarrow \infty.$$

In addition, one also has the convergence in distribution

$$\frac{3^{1/2}\{N_n^*(1) - (2n)^{1/2}\}}{(2n)^{1/4}} \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty.$$

However, such limiting results are only available for this special case, and our goal is to generalize such results to a broader class of weight distributions. We note that the available limiting results of the special case is obtained through the equivalence between the sequential monotone subsequence selection problem and the dynamic and stochastic knapsack problem with equal rewards (see Section 3.3). While we need a different proof strategy for the general case.

4.1.1 The Adaptive Heuristic $\hat{\pi}$ Revisited

Let us recall the dynamic and stochastic knapsack problem formulation of Chapter 2, and focus our attention to the case in which $p = r = 1$. That is, the items arrive deterministically in each period and have unitary rewards. The consumption function of (2.2) then becomes

$$\epsilon_k(x) = \sup \left\{ \epsilon \in [0, \infty) : \int_0^\epsilon w dF(w) \leq \frac{x}{k} \right\} \quad \text{for all } k \geq 1 \text{ and all } x \in [0, c]. \quad (4.3)$$

This consumption function in turn gives us the adaptive heuristic of Section 2.3, given as $\hat{\pi} = \{\hat{h}_n, \dots, \hat{h}_1\}$, where for all $k \in [n]$,

$$\hat{h}_k(x) = \min\{x, \epsilon_k(x)\} \quad \text{for all } x \in [0, c]. \quad (4.4)$$

Such adaptive heuristic uses time- and state-dependent thresholds for selection decisions. Specifically, for the t -th arriving item with weight realization W_t , if the current capacity level is x , then the policy $\hat{\pi}$ includes this item into the knapsack if and only if $W_t \leq \hat{h}_{n-t+1}(x)$. Thus, one can recursively define the stochastic process $\{X_t\}_{t \in [n]}$ of the remaining capacity under policy $\hat{\pi}$ by setting $X_0 = c$ and for all $t \in [n]$

$$X_t = \begin{cases} X_{t-1} & \text{if } W_t > \hat{h}_{n-t+1}(X_{t-1}) \\ X_{t-1} - W_t & \text{if } W_t \leq \hat{h}_{n-t+1}(X_{t-1}). \end{cases} \quad (4.5)$$

The number of selections $\hat{N}_n(c)$ under policy $\hat{\pi}$ is then given by

$$\hat{N}_n(c) = \sum_{t=1}^n \mathbb{1} \left\{ W_t \leq \hat{h}_{n-t+1}(X_{t-1}) \right\}.$$

We note that such definition of the number of selections is equivalent to the definition through stopping times in (4.2), and those stopping times given in (4.1) are the times at which selections occur.

Besides the remaining capacity process $\{X_t\}_{t \in [n]}$ and the number of selections $\hat{N}_n(c)$, the adaptive heuristic $\hat{\pi} = \{\hat{h}_n, \dots, \hat{h}_1\}$ has its associated value functions $\{\hat{v}_k(x)\}_{k \in [n]}$. With $\hat{v}_0(x) \equiv 0$, these value functions are recursively defined as

$$\hat{v}_k(x) = (1 - F(\hat{h}_k(x)))\hat{v}_{k-1}(x) + \int_0^{\hat{h}_k(x)} \{1 + \hat{v}_{k-1}(x - w)\} f(w) dw \quad \text{for all } k \in [n]. \quad (4.6)$$

The first-order performance of $\hat{\pi}$ is then given by the identity $\hat{v}_n(c) = \mathbb{E}[\hat{N}_n(c)]$. Lastly, by taking derivative of this last recursion, we obtain the recursion of the derivative of the value function

$$\hat{v}'_k(x) = \begin{cases} (1 - F(x))\hat{v}'_{k-1}(x) + \int_0^x \hat{v}'_{k-1}(x - w)f(w) dw + f(x)(1 - \hat{v}_{k-1}(x)) & \text{if } \epsilon_k(x) > x \\ (1 - F(\epsilon_k(x)))\hat{v}'_{k-1}(x) + \int_0^{\epsilon_k(x)} \hat{v}'_{k-1}(x - w)f(w) dw \\ + \frac{1}{k\epsilon_k(x)} \{1 + \hat{v}_{k-1}(x - \epsilon_k(x)) - \hat{v}_{k-1}(x)\} & \text{if } \epsilon_k(x) \leq x. \end{cases} \quad (4.7)$$

If the item-weight distribution F is continuous with density f , then one has the integral representation that for all $0 \leq x \leq k\mu \equiv k\mathbb{E}[W_1]$,

$$\int_0^{\epsilon_k(x)} wf(w) dw = \frac{x}{k}. \quad (4.8)$$

Two right-away observations are in order. First, for all $x \in [0, k\mu]$, one has that $0 \leq \frac{(k-1)x}{k} \leq (k-1)\mu$, so

$$\int_0^{\epsilon_{k-1}(\frac{(k-1)x}{k})} wf(w) dw = \frac{x}{k}.$$

Thus, we find that

$$\epsilon_{k-1}\left(\frac{(k-1)x}{k}\right) = \epsilon_k(x) \quad \text{for all } x \in [0, k\mu]. \quad (4.9)$$

Second, when F is continuous with density f , implicit function theorem tells us that $\epsilon_k(x)$ is differentiable on $(0, k\mu)$, and its derivative is given by

$$\epsilon'_k(x) = \frac{1}{k\epsilon_k(x)f(\epsilon_k(x))} \quad \text{for all } x \in (0, k\mu). \quad (4.10)$$

This last derivative representation then gives us the concavity of $F^\alpha(\epsilon_k(x))$ for any power $\alpha \in [0, 1]$. Such property is gathered in the following proposition.

Proposition 4.1 (Concavity). *Let $0 \leq \alpha \leq 1$ be fixed, then for any $k \geq 1$ and any $x \in [0, k\mu]$, the map $x \mapsto F^\alpha(\epsilon_k(x))$ is concave in x .*

Proof. Recall from (4.10) that for any $x \in [0, k\mu]$, we have

$$\epsilon'_k(x) = \frac{1}{k\epsilon_k(x)f(\epsilon_k(x))}.$$

As a result, for all $0 \leq \alpha \leq 1$ fixed, the derivative of $F^\alpha(\epsilon_k(x))$ is

$$\frac{d}{dx} F^\alpha(\epsilon_k(x)) = \frac{\alpha}{k\epsilon_k(x)F^{1-\alpha}(\epsilon_k(x))},$$

which is monotone decreasing in x . Therefore, the map $x \mapsto F^\alpha(\epsilon_k(x))$ is concave in x . \square

We close this section with the martingale structure of a normalized remaining capacity process induced by $\hat{\pi}$, as well as an easy implication of such martingale structure. Specifically, if we consider the stopping time

$$\eta \equiv \min \{t \in \{0, 1, \dots, n\} : \epsilon_{n-t}(X_t) > X_t \text{ or } X_t > (n-t)\mu\}, \quad (4.11)$$

then with appropriate normalization, the remaining capacity process up until η forms a martingale. Such observation is made precise in the following proposition.

Proposition 4.2 (Martingale structure associated with the remaining capacity process). *Let the stopping time η be defined as in (4.11), then the normalized remaining capacity process up until η , i.e., $\{(n-t)^{-1}X_t : t = 0, 1, \dots, \eta\}$ is a martingale. That is,*

$$\mathbb{E} \left[\frac{X_t}{n-t} \middle| \mathcal{F}_{t-1} \right] = \frac{X_{t-1}}{n-t+1} \quad \text{for all } t = 1, \dots, \eta.$$

Proof. For all $t \in [n]$, the construction of $\hat{\pi}$ gives us that

$$\mathbb{E}[X_t | \mathcal{F}_{t-1}] = X_{t-1} - \int_0^{\hat{h}_{n-t+1}(X_{t-1})} wf(w) dw.$$

When $t \leq \eta$ in particular, the definition (4.11) tells us two more conditions. First, for all $1 \leq t \leq \eta$,

$$\hat{h}_{n-t+1}(X_{t-1}) = \min\{\epsilon_{n-t+1}(X_{t-1}), X_{t-1}\} = \epsilon_{n-t+1}(X_{t-1}).$$

Second, for all $1 \leq t \leq \eta$, the remaining capacity always satisfies $X_{t-1} \leq (n-t+1)\mu$. This, together with (4.8), implies that

$$\int_0^{\epsilon_{n-t+1}(X_{t-1})} wf(w) dw = \frac{X_{t-1}}{n-t+1}.$$

Putting the last three identities altogether, we obtain that for all $1 \leq t \leq \eta$,

$$\mathbb{E}[X_t \mid \mathcal{F}_{t-1}] = \frac{(n-t)X_{t-1}}{n-t+1},$$

or equivalently,

$$\mathbb{E}\left[\frac{X_t}{n-t} \mid \mathcal{F}_{t-1}\right] = \frac{X_{t-1}}{n-t+1}.$$

Therefore, $\{(n-t)^{-1}X_t : t = 0, 1, \dots, \eta\}$ is a martingale. \square

This last martingale structure, together with Proposition 4.1, gives us the following series upper bound.

Proposition 4.3 (Series upper bound). *Let the remaining capacity process $\{X_t : t = 0, 1, \dots, n\}$ under the adaptive heuristic $\hat{\pi}$ be defined as in (4.5), and let the stopping time η be defined as in (4.11). Then for any fixed $0 < \alpha \leq 1$, we have that for all $0 \leq t \leq \eta - 1$,*

$$\sum_{j=t+1}^{\eta} \mathbb{E}\left[\frac{F^{\alpha}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{1-\alpha}} \mid \mathcal{F}_t\right] \leq \alpha^{-1} \{(n-t)F(\epsilon_{n-t}(X_t))\}^{\alpha} \quad a.s.$$

Proof. From Proposition 4.1 we know that $F^{\alpha}(\epsilon_k(x))$ is concave in x , so we can apply Jensen's inequality to obtain that for $j \geq 2$,

$$\mathbb{E}[F^{\alpha}(\epsilon_{n-j+1}(X_{j-1})) \mid \mathcal{F}_{j-2}] \leq F^{\alpha}(\epsilon_{n-j+1}(\mathbb{E}[X_{j-1} \mid \mathcal{F}_{j-2}])).$$

For any realization of the stopping time η , Proposition 4.2 tells us that for $2 \leq j \leq \eta$,

$$\mathbb{E}\left[\frac{X_{j-1}}{n-j+1} \mid \mathcal{F}_{j-2}\right] = \frac{X_{j-2}}{n-j+2},$$

which in turn implies that

$$\mathbb{E}[F^{\alpha}(\epsilon_{n-j+1}(X_{j-1})) \mid \mathcal{F}_{j-2}] \leq F^{\alpha}\left(\epsilon_{n-j+1}\left(\frac{(n-j+1)X_{j-2}}{n-j+2}\right)\right) = F^{\alpha}(\epsilon_{n-j+2}(X_{j-2})),$$

where in the last equality we have used identity (4.9). With repeated applications of Jensen's inequality and the martingale structure, we finally obtain that for each summand of the series,

$$\mathbb{E}[F^{\alpha}(\epsilon_{n-j+1}(X_{j-1})) \mid \mathcal{F}_i] \leq F^{\alpha}(\epsilon_{n-i}(X_i)).$$

Therefore, the series is bounded from above by

$$\sum_{j=t+1}^{\eta} \mathbb{E} \left[\frac{F^{\alpha}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{1-\alpha}} \middle| \mathcal{F}_t \right] \leq F^{\alpha}(\epsilon_{n-t}(X_t)) \sum_{j=t+1}^{\tau} \frac{1}{(n-t+1)^{1-\alpha}}.$$

Applying the crude bound $\eta \leq n$ as well as the basic series inequality $\sum_{j=t+1}^n (n-j+1)^{\alpha-1} \leq \alpha^{-1}(n-t)^{\alpha}$, we obtain that

$$\sum_{j=t+1}^{\eta} \mathbb{E} \left[\frac{F^{\alpha}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{1-\alpha}} \middle| \mathcal{F}_t \right] \leq \alpha^{-1} \{(n-t)F(\epsilon_{n-t}(X_t))\}^{\alpha}.$$

Since the above argument holds for any realization of η , it holds almost surely. \square

4.1.2 Main Result

When it comes to performance analysis of the dynamic and stochastic knapsack problems, regularity conditions are common practice in the literature. Since we are interested in the limiting regime in which the knapsack capacity is fixed while the number of arriving items grows, the behavior of the item-weight distribution F around zero is particularly important to the asymptotic analysis. We specify the regularity conditions for our analysis in the following definition.

Definition 4.1 (Regular distributions.). Any non-negative continuous distribution F with density f is said to be *regular* if it satisfies the following conditions on $(0, \bar{w})$ for some $\bar{w} > 0$.

(i) BOUNDED DENSITY.

$$0 < m_{\ell} \leq f(w) \leq m_u < \infty \quad \text{for all } w \in (0, \bar{w}). \quad (4.12)$$

(ii) SMOOTHNESS AND EXISTENCE OF LIMIT. The density f is differentiable and there exists constant $r \in (0, \infty)$ such that

$$\lim_{w \rightarrow 0^+} \frac{wf'(w)}{f(w)} = \frac{1}{r} - 1. \quad (4.13)$$

(iii) CONVERGENCE RATE. In the neighborhood of zero,

$$\left| \frac{F(w)}{wf(w)} - r \right| \leq Mw^{3/4}, \quad \text{for all } w \in (0, \bar{w}). \quad (4.14)$$

Logarithmic Regret Bound

As a quick implication of Condition (ii), we have a lower bound on the first-order performance of the adaptive heuristic $\hat{\pi}$. Specifically, from Condition (ii), one has that

$$(wf(w))' = f(w) \left(\frac{wf'(w)}{f(w)} + 1 \right) \geq 0$$

in a neighborhood of zero, hence the function $wf(w)$ is monotone increasing. Moreover, by L'Hospital's rule and (4.13), we know that

$$\lim_{w \rightarrow 0^+} \frac{wF(w)}{\int_0^w F(y) dy} = \lim_{w \rightarrow 0^+} \frac{F(w) + wf(w)}{F(w)} = \frac{1+r}{r} > 1.$$

Then it follows from Theorem 2.1 and Lemma 2.2 that

$$nF(\epsilon_n(c)) - \mathbb{E} \left[\hat{N}_n(c) \right] = O(\log n). \quad (4.15)$$

We state our main result in the following theorem.

Theorem 4.1 (Main result). *Consider the dynamic and stochastic knapsack problem with unitary rewards, initial capacity c and regular item-weight distribution F as defined in Definition 4.1. Let the associated constant r be as in (4.13), if we set*

$$\gamma(r) \equiv \frac{2r^2}{(1+r)(1+2r)},$$

then the number of selections $\hat{N}_n(c)$ of the adaptive heuristic policy $\hat{\pi}$ satisfies the convergence in distribution

$$\frac{\hat{N}_n(c) - nF(\epsilon_n(c))}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty.$$

Moreover, if the distribution F satisfies the growth condition

$$\frac{nF(\epsilon_n(c))}{(\log n)^2} \rightarrow \infty \quad \text{as } n \rightarrow \infty, \quad (4.16)$$

then the number of selections $N_n^(c)$ under the optimal dynamic programming policy π^* satisfies the same convergence in distribution*

$$\frac{N_n^*(c) - nF(\epsilon_n(c))}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty.$$

The item-weight distribution F plays a key role in the proof of our main result. Next, we collect the properties of F when it satisfies the regularity conditions in Definition 4.1.

4.1.3 Properties of the Regular Class

Definition 4.1 imposes regularity conditions on the item-weight distribution F . Such regularity conditions have several implications that are easy to use in the analysis. We summarize the implications below. The first implication is that of Condition (ii) in Definition 4.1.

Proposition 4.4 (Implication of Condition (ii), behavior near zero). *Condition (ii) in Definition 4.1 of the regular class implies that there exists an interval $(0, \tilde{w})$ such that for all $w \in (0, \tilde{w})$, the two auxiliary functions*

$$T_1(w) \equiv wf(w), \quad T_2(w) \equiv \frac{\left(\int_0^w tf(t) dt\right)^2}{w^2 f(w)} \quad (4.17)$$

are monotone increasing; and

$$\max \left\{ \frac{wf'(w)}{f(w)}, \frac{wf(w)}{F(w)} \right\} \leq M \quad (4.18)$$

for some constant M that depends only on F .

Proof. From (4.13) we know that there exists an interval $(0, \bar{w}_1)$ such that for all $w \in (0, \bar{w}_1)$,

$$\frac{wf'(w)}{f(w)} + 1 \geq 0,$$

which in turn implies that $T_1'(w) = wf'(w) + f(w) \geq 0$ and hence $T_1(w)$ is increasing in $w \in (0, \bar{w}_1)$.

On the other hand, the limit (4.13) also implies the following limit

$$\lim_{w \rightarrow 0^+} \left\{ \frac{w^2 f(w)}{\int_0^w tf(t) dt} - \frac{wf'(w)}{2f(w)} - 1 \right\} = \lim_{w \rightarrow 0^+} \frac{2wf(w) + w^2 f'(w)}{wf(w)} - a/2 - 1 = 1 + a/2 > 0.$$

As a result, there exists an interval $(0, \bar{w}_2)$ such that for any $w \in (0, \bar{w}_2)$,

$$\frac{w^2 f(w)}{\int_0^w tf(t) dt} - \frac{wf'(w)}{2f(w)} - 1 \geq 0.$$

Now if we compute the derivative of $T_2(w)$ in $w \in (0, \bar{w}_2)$, we would get

$$\begin{aligned} T_2'(w) &= \frac{2 \int_0^w t f(t) dt \cdot w^3 f^2(w) - \left(\int_0^w t f(t) dt \right)^2 (w^2 f'(w) + 2w f(w))}{w^4 f^2(w)} \\ &= \frac{2 \left(\int_0^w t f(t) dt \right)^2}{w^3 f(w)} \left(\frac{w^2 f(w)}{\int_0^w t f(t) dt} - \frac{w f'(w)}{2f(w)} - 1 \right) \geq 0, \end{aligned}$$

so $T_2(w)$ is increasing in $w \in (0, \bar{w}_2)$.

As for the upper bound (4.18), if we set $\tilde{w} = \min\{\bar{w}_1, \bar{w}_2\}$ and consider the continuous function $w f'(w)/f(w)$, then the existence of the limit (4.13) tells us that this continuous function defined on $w \in (0, \tilde{w}]$ can be extended to the compact interval $w \in [0, \tilde{w}]$. Since a continuous function on a compact interval is always bounded, we have that

$$\sup_{w \in [0, \tilde{w}]} \frac{w f'(w)}{f(w)} = M < \infty,$$

and apparently the constant M depends only on F . The same argument applies to the other continuous function $w f(w)/F(w)$ as well. In fact, the existence of the following finite limit

$$\lim_{w \rightarrow 0+} \frac{w f(w)}{F(w)} = \lim_{w \rightarrow 0+} \frac{w f'(w) + f(w)}{f(w)} = a + 1 < \infty$$

implies that the domain of $\frac{w f(w)}{F(w)}$ can also be extended to the compact interval $w \in [0, \tilde{w}]$. Hence, we have that

$$\sup_{w \in [0, \tilde{w}]} \frac{w f(w)}{F(w)} = M < \infty,$$

with M depending only on F . We now have all the claimed properties of regular distribution within $(0, \tilde{w})$. \square

The second useful implication is that of Condition (iii) in Definition 4.1.

Proposition 4.5 (Implication of Condition (iii)). *Condition (iii) in the Definition 4.1 implies that there exist positive constants M and \bar{w} such that*

$$\max \left\{ \left| \frac{\int_0^w F(y) dy}{w F(w)} - \frac{r}{1+r} \right|, \left| \frac{\int_0^w y F(y) dy}{w^2 F(w)} - \frac{r}{1+2r} \right| \right\} \leq M w^{3/4} \quad \text{for all } w \in (0, \bar{w}). \quad (4.19)$$

Proof. To prove (4.19), we start with the first maximand of the left-hand side. From equation (4.14) we know that for all $w \in (0, \bar{w})$,

$$\left| \frac{\int_0^w F(y) dy}{\int_0^w yf(y) dy} - r \right| \leq Mw^{3/4}.$$

If we multiply both sides with $(1+r)^{-1} \int_0^w yf(y) dy$, we obtain that

$$(1+r)^{-1} \left| \int_0^w F(y) dy - r \int_0^w yf(y) dy \right| \leq (1+r)^{-1} Mw^{3/4} \int_0^w yf(y) dy.$$

Note that with an integration by parts, the left-hand side of this last inequality can be re-written as

$$\begin{aligned} (1+r)^{-1} \left| \int_0^w F(y) dy - r \int_0^w yf(y) dy \right| &= (1+r)^{-1} \left| \int_0^w F(y) dy - r \left(wF(w) - \int_0^w F(y) dy \right) \right| \\ &= \left| \int_0^w F(y) dy - \frac{r}{1+r} wF(w) \right|, \end{aligned}$$

so we have

$$\left| \int_0^w F(y) dy - \frac{r}{1+r} wF(w) \right| \leq (1+r)^{-1} Mw^{3/4} \int_0^w yf(y) dy.$$

By dividing both sides with $wF(w)$, we finally obtain that

$$\left| \frac{\int_0^w F(y) dy}{wF(w)} - \frac{r}{1+r} \right| \leq \frac{\int_0^w yf(y) dy}{(1+r)wF(w)} Mw^{3/4} \leq Mw^{3/4},$$

where the last inequality is because $r > 0$ and $\int_0^w yf(y) dy \leq wF(w)$.

As for the second maximand in (4.19), we know from (4.14) that

$$\left| \frac{wF(w)}{w^2 f(w)} - r \right| \leq Mw^{3/4} \quad \text{for all } w \in (0, \bar{w}),$$

so we further have that for all $w \in (0, \bar{w})$,

$$\left| \frac{\int_0^w yF(y) dy}{\int_0^w y^2 f(y) dy} - r \right| \leq Mw^{3/4}.$$

Then we multiply both sides with $(1+2r)^{-1} \int_0^w y^2 f(y) dy$ and use integration by parts to the integral $\int_0^w y^2 f(y) dy$ on the left-hand side to obtain that

$$\left| \int_0^w yF(y) dy - \frac{r}{1+2r} w^2 F(w) \right| \leq (1+2r)^{-1} Mw^{3/4} \int_0^w y^2 f(y) dy.$$

Finally, by dividing both sides with $w^2 F(w)$, we obtain that

$$\left| \frac{\int_0^w y F(y) dy}{w^2 F(w)} - \frac{r}{1+2r} \right| \leq \frac{\int_0^w y^2 f(y) dy}{(1+2r)w^2 F(w)} M w^{3/4} \leq M w^{3/4},$$

completing the second half of (4.19). \square

To take full use of the properties of the regular class, we define constant

$$K_0 \equiv \max \left\{ \left\lceil \frac{c}{\int_0^{\min\{\bar{w}, \tilde{w}\}} w f(w) dw} \right\rceil, 3 \right\}, \quad (4.20)$$

and restrict ourselves into the problem instances with more than K_0 items. In fact, the choice of K_0 guarantees us that the consumption function $\epsilon_k(x)$ always satisfies $0 \leq \epsilon_k(x) \leq \min\{\bar{w}, \tilde{w}\}$ for all $x \in [0, c]$ whenever $k \geq K_0$.

The smoothness assumption of the item-weight distribution F implies quite rich analytical properties, which are heavily used in the sequel analysis. We close this section with these properties of regular class distributions.

Proposition 4.6 (Properties of regular class.). *Let F be a given distribution that belongs to the regular class, then there exists constant M_F such that for all $x \in (0, c]$ and all $k \geq K_0$, the following properties hold.*

(i) BOUNDED RATIOS.

$$1 \leq \frac{k \epsilon_k(x) F(\epsilon_k(x))}{x} \leq M_F, \quad (4.21)$$

and

$$1 \leq \frac{k \epsilon_k^2(x) f(\epsilon_k(x))}{x} \leq M_F. \quad (4.22)$$

(ii) DIFFERENCE IN PROPHET UPPER BOUND. *For any $w \in [0, x]$,*

$$0 \leq k F(\epsilon_k(x)) - k F(\epsilon_k(x-w)) - \frac{w}{\epsilon_k(x)} \leq \frac{k F(\epsilon_k(x)) w^2}{x^2}. \quad (4.23)$$

(iii) A MONOTONE RELATION.

$$k^2 \epsilon_k^2(x) f(\epsilon_k(x)) \leq (k+1)^2 \epsilon_{k+1}^2(x) f(\epsilon_{k+1}(x)). \quad (4.24)$$

(iv) TAYLOR-EXPANSION TYPE BOUNDS. For any $w \in [0, x/2]$,

$$\frac{\epsilon_k(x)}{\epsilon_k(x-w)} - 1 \leq \frac{w}{k\epsilon_k^2(x)f(\epsilon_k(x))} + M_F \frac{w^2}{x^2}. \quad (4.25)$$

Moreover, for any fixed $\alpha \in (0, 1/2)$, there exists a constant $M_{F,\alpha}$ such that

$$\frac{x^{1-\alpha}\epsilon_k^\alpha(x)}{(x-w)^{1-\alpha}\epsilon_k^\alpha(x-w)} - 1 \leq \left(\frac{1-\alpha}{x} + \frac{\alpha}{k\epsilon_k^2(x)f(\epsilon_k(x))} \right) w + M_{F,\alpha} \frac{w^2}{x^2}. \quad (4.26)$$

Proof. (i) The upper bound in (4.21) was proved in Lemma 2.2, while the lower bound in (4.21) is an easy consequence of (4.8)—if we apply integration by parts in (4.8) to get

$$\epsilon_k(x)F(\epsilon_k(x)) - \int_0^{\epsilon_k(x)} F(w) dw = \frac{x}{k},$$

and drop the positive integral, we obtain the inequality

$$\epsilon_k(x)F(\epsilon_k(x)) \geq \frac{x}{k},$$

which in turn implies the lower bound in (4.21).

As for the second inequality (4.22), we know that when $k \geq K_0$, the consumption function satisfies that $\epsilon_k(x) \leq \tilde{w}$ for all $x \in [0, c]$. Therefore, the monotonicity of $wf(w)$ from (4.17) implies that

$$\frac{x}{k} = \int_0^{\epsilon_k(x)} wf(w) dw \leq \int_0^{\epsilon_k(x)} \epsilon_k(x)f(\epsilon_k(x)) dw = \epsilon_k^2(x)f(\epsilon_k(x)),$$

which is the left inequality of (4.22). On the other hand, with some algebra, we know that

$$\frac{k\epsilon_k^2(x)f(\epsilon_k(x))}{x} = \frac{k\epsilon_k(x)F(\epsilon_k(x))}{x} \frac{\epsilon_k(x)f(\epsilon_k(x))}{F(\epsilon_k(x))}.$$

Combining the first bound (4.21) and the bound (4.18) from Proposition 4.4, we get the right inequality of the second bound

$$\frac{k\epsilon_k^2(x)f(\epsilon_k(x))}{x} \leq M_F.$$

(ii) This was proved in Equation (2.28) of Proposition 2.2.

(iii) By the choice of K_0 and the definition of the consumption function, we have that $\epsilon_{k+1}(x) \leq \epsilon_k(x) \leq \tilde{w}$ for all $x \in [0, c]$. Then the monotone relation (4.24) is just an easy corollary of the

monotonicity of the function $T_2(w)$ in (4.17) from Proposition 4.4. In fact, since $T_2(w)$ is increasing in $w \in (0, \tilde{w}]$, we have that $T_2(\epsilon_{k+1}(x)) \leq T_2(\epsilon_k(x))$, which is

$$\frac{x/(k+1)}{(k+1)\epsilon_{k+1}^2(x)f(\epsilon_{k+1}(x))} \leq \frac{x/k}{k\epsilon_k^2(x)f(\epsilon_k(x))}.$$

This last inequality in turn implies that

$$k^2\epsilon_k^2(x)f(\epsilon_k(x)) \leq (k+1)^2\epsilon_{k+1}^2(x)f(\epsilon_{k+1}(x)),$$

as desired.

(iv) We first prove the first inequality (4.25), and the second inequality (4.26) follows a similar argument. The proof is basically Taylor expansion together with a verification of the boundedness of the second order derivative. For any $w \in [0, x/2]$, we first show that $\epsilon_k(x)/\epsilon_k(x-w) \leq 2$. In fact, $0 \leq w \leq x/2$ implies that

$$1 \leq \frac{x}{x-w} \leq 2. \quad (4.27)$$

From (4.8), we also know that

$$\int_{\epsilon_k(x/2)}^{\epsilon_k(x)} wf(w) dw = \frac{x}{2k} = \int_0^{\epsilon_k(x/2)} wf(w) dw. \quad (4.28)$$

The monotonicity of $wf(w)$ from Proposition 4.4 implies that

$$\begin{aligned} \int_{\epsilon_k(x/2)}^{\epsilon_k(x)} wf(w) dw &\geq \epsilon_k(x/2)f(\epsilon_k(x/2))\{\epsilon_k(x) - \epsilon_k(x/2)\}, \\ \int_0^{\epsilon_k(x/2)} wf(w) dw &\leq \epsilon_k(x/2)f(\epsilon_k(x/2))\epsilon_k(x/2). \end{aligned}$$

By connecting these two inequalities through (4.28) we obtain that

$$\epsilon_k(x) \leq 2\epsilon_k(x/2),$$

which in turn tells us that for all $w \in [0, x/2]$,

$$\frac{\epsilon_k(x)}{\epsilon_k(x-w)} \leq \frac{\epsilon_k(x)}{\epsilon_k(x/2)} \leq 2. \quad (4.29)$$

Next, we calculate the first two derivatives of $1/\epsilon_k(x)$, and, with some algebra, obtain that

$$\begin{aligned} \frac{d}{dx} \frac{1}{\epsilon_k(x)} &= -\frac{1}{k\epsilon_k^3(x)f(\epsilon_k(x))}, \\ \frac{d^2}{dx^2} \frac{1}{\epsilon_k(x)} &= \frac{1}{k^2\epsilon_k^5(x)f^2(\epsilon_k(x))} \left(3 + \frac{\epsilon_k(x)f'(\epsilon_k(x))}{f(\epsilon_k(x))} \right). \end{aligned}$$

Therefore for any $w \in [0, x/2]$, the second derivative satisfies that

$$\begin{aligned} & x^2 \epsilon_k^2(x) \cdot \frac{d^2}{dw^2} \frac{1}{\epsilon_k(x-w)} \\ &= \frac{x^2 \epsilon_k(x)}{(x-w)^2 \epsilon_k(x-w)} \frac{(x-w)^2}{k^2 \epsilon_k^4(x-w) f^2(\epsilon_k(x-w))} \left(3 + \frac{\epsilon_k(x-w) f'(\epsilon_k(x-w))}{f(\epsilon_k(x-w))} \right) \leq M_F, \end{aligned}$$

where we have used (4.22) for the boundedness of the second factor; (4.18) for the boundedness of the last factor; and the fact

$$\frac{x^2 \epsilon_k(x)}{(x-w)^2 \epsilon_k(x-w)} \leq 8$$

that we have just proved (a combination of (4.27) and (4.29)). As a result, the second order Taylor expansion upper bound for $1/\epsilon_k(x-w)$ at $w=0$ reads

$$\frac{1}{\epsilon_k(x-w)} \leq \frac{1}{\epsilon_k(x)} + \frac{w}{k \epsilon_k^3(x) f(\epsilon_k(x))} + \frac{M_F w^2}{x^2 \epsilon_k(x)},$$

which is equivalent to inequality (4.25) as we wanted.

On the other hand, we can expand the other function $\{(x-w)^{\alpha-1} \epsilon_k^{-\alpha}(x-w)\}$ at $w=0$ as well.

With some algebra and simplifications, we obtain the first two derivatives

$$\begin{aligned} \frac{d}{dx} ((x-w)^{\alpha-1} \epsilon_k^{-\alpha}(x-w)) &= (\alpha-1) x^{\alpha-2} \epsilon_k^{-\alpha}(x) - \alpha k^{-1} x^{\alpha-1} \epsilon_k^{-\alpha-2}(x) f^{-1}(\epsilon_k(x)), \\ \frac{d^2}{dx^2} ((x-w)^{\alpha-1} \epsilon_k^{-\alpha}(x-w)) &= (1-\alpha)(2-\alpha) x^{\alpha-3} \epsilon_k^{-\alpha}(x) + 2\alpha(1-\alpha) k^{-1} x^{\alpha-2} \epsilon_k^{-\alpha-2}(x) f^{-1}(\epsilon_k(x)) \\ &\quad + \alpha(\alpha+2) k^{-2} x^{\alpha-1} \epsilon_k^{-\alpha-4}(x) f^{-2}(\epsilon_k(x)) + \alpha k^{-2} x^{\alpha-1} \epsilon_k^{-\alpha-3}(x) f'(\epsilon_k(x)) f^{-2}(\epsilon_k(x)). \end{aligned}$$

Similar to what we did for the previous expansion, it is simple but tedious work to check that for any $w \in [0, x/2]$,

$$x^{3-\alpha} \epsilon_k^\alpha(x) \left\{ \frac{d^2}{dw^2} ((x-w)^{\alpha-1} \epsilon_k^{-\alpha}(x-w)) \right\} \leq M_{F,\alpha},$$

hence the second order Taylor expansion bound holds, that is, for any $w \in [0, x/w]$,

$$\begin{aligned} (x-w)^{\alpha-1} \epsilon_k^{-\alpha}(x-w) &\leq x^{\alpha-1} \epsilon_k^{-\alpha}(x) + \{(1-\alpha) x^{\alpha-2} \epsilon_k^{-\alpha}(x) + \alpha k^{-1} x^{\alpha-1} \epsilon_k^{-\alpha-2}(x) f^{-1}(\epsilon_k(x))\} w \\ &\quad + \frac{M_{F,\alpha} w^2}{x^{11/4} \epsilon_k^{1/4}(x)}, \end{aligned}$$

which is equivalent to the second inequality (4.26) as we wanted. \square

4.2 Non-asymptotic Derivative Bounds

The analysis of the value function derivative $\hat{v}'_k(x)$ plays a crucial role in our proof of Theorem 4.1. In a nutshell, if the derivative $\hat{v}'_k(x)$ can be tightly estimated, then the fluctuations along the sequential decision process can also be well controlled. Such important intermediate step is given in the following lemma.

Lemma 4.1 (Non-asymptotic derivative bounds.). *If the item-weight distribution F belongs to the regular class in Definition 4.1, then there exists a constant m that depends only on F such that for all $k \geq K_0$, the derivative $\hat{v}'_k(x)$ is bounded by*

$$\ell_k(m, x) \leq \hat{v}'_{k-1}(x) \leq u_k(m, x) \quad \text{for all } x \in [0, c],$$

where the two functions are given by

$$u_k(m, x) = \frac{1}{\epsilon_k(x)} + \frac{m}{x^{3/4}\epsilon_k^{1/4}(x)}, \quad \text{and} \quad \ell_k(m, x) = \frac{1}{\epsilon_k(x)} - \frac{m}{x^{3/4}\epsilon_k^{1/4}(x)}.$$

The proof of Lemma 4.1 is based on an induction argument that consists of two parts. Essentially, we split up the domain $x \in [0, c]$ into two regions according to the ratio $\frac{x}{\epsilon_k(x)}$ and analyze the two regions separately. Specifically, we define two constants

$$M_0 \equiv \max \left\{ 16M_F, \frac{m_u}{2m_\ell}, 2 \right\} \tag{4.30}$$

$$m_0 \equiv \max \left\{ \frac{8(M_F + 1)}{M_0^{1/4}}, \frac{m_u(M_F M_0 + 1) [M_0(1 + K_0^{-1})]^{7/4}}{m_\ell}, \right. \\ \left. [M_0(1 + K_0^{-1})]^{3/4} + \frac{m_u(M_F M_0 + 1) [M_0(1 + K_0^{-1})]^{7/4}}{m_\ell}, 1 \right\}, \tag{4.31}$$

and consider for $k \geq K_0 + 1$ the following two regions

$$\Omega_k \equiv \{x \in [0, c] : x \leq M_0 \epsilon_k(x)\}, \quad \text{and} \quad \Omega_k^c \equiv \{x \in [0, c] : x > M_0 \epsilon_k(x)\}.$$

The first part of the induction argument for proving Lemma 4.1 is to control the derivative $\hat{v}'_k(x)$ in the region Ω_k . We have the following proposition.

Proposition 4.7 (Derivative bounds when capacity is small). *If the item weight distribution F belongs to the regular class, then whenever $k \geq K_0 + 1$ and $x \leq M_0 \epsilon_{k-1}(x)$, we have the derivative*

bound

$$\ell_k(m, x) \leq \hat{v}'_{k-1}(x) \leq u_k(m, x) \quad \text{for all } x \in [0, c],$$

for all $m \geq m_0$ defined in (4.31).

Proof. Step 1. We start with the recursion of $\hat{v}'_k(x)$ and prove by induction that

$$\max_{x \in \Omega_k} |\hat{v}'_k(x)| \leq k(M_F M_0 + 1)m_u/2 \text{ for all } k \geq 1,$$

where $\Omega_k = \{x \in [0, c] : x \leq M_0 \epsilon_k(x)\}$.

Before the induction proof, let us prepare two estimations that will be used later. First, from the integral representation of the consumption function (4.8) and the boundedness of the density f from (4.12), we know that

$$\frac{x}{k} = \int_0^{\epsilon_k(x)} w f(w) dw \leq m_u \int_0^{\epsilon_k(x)} w dw = \frac{m_u \epsilon_k^2(x)}{2},$$

which yields the inequality

$$\frac{1}{k \epsilon_k(x)} \leq \frac{m_u \epsilon_k(x)}{2x}. \quad (4.32)$$

The second estimation is that when $x \leq M_0 \epsilon_{k-1}(x)$, from Theorem 2.1 we obtain that for all $y \in [0, x]$,

$$0 \leq \hat{v}_{k-1}(y) \leq (k-1)F(\epsilon_{k-1}(y)) \leq (k-1)F(\epsilon_{k-1}(x)) = \frac{(k-1)\epsilon_{k-1}(x)F(\epsilon_{k-1}(x))}{x} \frac{x}{\epsilon_{k-1}(x)}.$$

From (4.21) we know that the first factor of this last right-hand side is bounded by M_F ; and the second factor is bounded by M_0 since $x \leq M_0 \epsilon_{k-1}(x)$. Hence we have that for all x such that $x \leq M_0 \epsilon_{k-1}(x)$,

$$0 \leq \hat{v}_{k-1}(y) \leq M_F M_0 \quad \text{for all } y \in [0, x]. \quad (4.33)$$

Now we are in the position to prove by induction the upper bound $|\hat{v}'_k(x)| \leq k(M_F M_0 + 1)m_u/2$ for all $x \in \Omega_k = \{x \in [0, c] : x \leq M_0 \epsilon_k(x)\}$. Recall the recursion of $\hat{v}'_k(x)$ in (4.7). Since it has two different representations depending on the relation between $\epsilon_k(x)$ and x , we further divide the region $\Omega_{k-1} \equiv \{x \in [0, c] : x \leq M_0 \epsilon_{k-1}(x)\}$ into two sub-regions: $0 \leq x < \epsilon_k(x)$ and $\epsilon_k(x) \leq x \leq M_0 \epsilon_{k-1}(x)$.

On the one hand, if $0 \leq x < \epsilon_k(x)$, then we take absolute value on both sides of the first line of (4.7) and use triangle inequality to obtain that

$$\begin{aligned} |\hat{v}'_k(x)| &\leq (1 - F(x)) |\hat{v}'_{k-1}(x)| + \int_0^x |\hat{v}'_{k-1}(x-w)| f(w) dw + f(x) |1 - \hat{v}_{k-1}(x)| \\ &\leq \max_{y \in [0, x]} |\hat{v}'_{k-1}(y)| + f(x) |1 - \hat{v}_{k-1}(x)|. \end{aligned}$$

From (4.21) we know that

$$\hat{v}_{k-1}(x) \leq (k-1)F(\epsilon_{k-1}(x)) \leq \frac{M_F x}{\epsilon_{k-1}(x)} \leq \frac{M_F x}{\epsilon_k(x)} \leq M_F,$$

so if we apply this last estimation on $\hat{v}_{k-1}(x)$ and (4.12) to the inequality of $|\hat{v}'_k(x)|$, we find the upper bound

$$|\hat{v}'_k(x)| \leq \max_{y \in [0, x]} |\hat{v}'_{k-1}(y)| + f(x) |1 - \hat{v}_{k-1}(x)| \leq \max_{y \in [0, x]} |\hat{v}'_{k-1}(y)| + M_F m_u.$$

On the other hand, if $\epsilon_k(x) \leq x \leq M_0 \epsilon_{k-1}(x)$, then we take absolute value on both sides of the second line of (4.7) and use triangle inequality to obtain that

$$\begin{aligned} |\hat{v}'_k(x)| &\leq (1 - F(\epsilon_k(x))) |\hat{v}'_{k-1}(x)| + \int_0^{\epsilon_k(x)} |\hat{v}'_{k-1}(x-w)| f(w) dw \\ &\quad + \frac{1}{k\epsilon_k(x)} |1 + \hat{v}_{k-1}(x - \epsilon_k(x)) - \hat{v}_{k-1}(x)| \\ &\leq \max_{y \in [x - \epsilon_k(x), x]} |\hat{v}'_{k-1}(y)| + \frac{1}{k\epsilon_k(x)} |1 + \hat{v}_{k-1}(x - \epsilon_k(x)) - \hat{v}_{k-1}(x)|. \end{aligned}$$

Next, we apply (4.32) and (4.33) to this last right hand to obtain the upper bound

$$|\hat{v}'_k(x)| \leq \max_{y \in [x - \epsilon_k(x), x]} |\hat{v}'_{k-1}(y)| + (M_F M_0 + 1)m_u/2.$$

Combining these last two inequalities on $|\hat{v}'_k(x)|$, we obtain that for all $x \in \Omega_{k-1}$,

$$\begin{aligned} |\hat{v}'_k(x)| &\leq \max_{y \in [x - \hat{h}_k(x), x]} |\hat{v}'_{k-1}(y)| + \max\{M_F m_u, (M_F M_0 + 1)m_u/2\} \\ &= \max_{y \in [x - \hat{h}_k(x), x]} |\hat{v}'_{k-1}(y)| + (M_F M_0 + 1)m_u/2. \end{aligned} \tag{4.34}$$

Lastly, we define $x_k \equiv \sup\{x \in [0, c] : x \in \Omega_k\}$ and prove by induction that for all $k \geq 1$,

$$\max_{x \in \Omega_k} |\hat{v}'_k(x)| \leq \max_{x \in [0, x_k]} |\hat{v}'_k(x)| \leq k(M_F M_0 + 1)m_u/2.$$

In fact, the nesting relation $\Omega_k \subset \Omega_{k-1}$ for all $k \geq 2$ tells us that $x_k \leq x_{k-1}$ for all $k \geq 2$. Since $\hat{v}_1(x) = F(x)$ for all $x \in [0, c]$, we have from (4.12) that

$$\max_{x \in [0, x_1]} |\hat{v}'_1(x)| \leq \max_{x \in [0, c]} |\hat{v}'_1(x)| = \max_{x \in [0, c]} f(x) \leq m_u \leq (M_F M_0 + 1)m_u/2,$$

so the induction base is satisfied. Then, if $\max_{x \in [0, x_{k-1}]} |\hat{v}'_{k-1}(x)| \leq (k-1)(M_F M_0 + 1)m_u/2$, we have from (4.34) that

$$\begin{aligned} \max_{x \in [0, x_k]} |\hat{v}'_k(x)| &\leq \max_{x \in [0, x_k]} \max_{y \in [\hat{h}_k(x), x]} |\hat{v}'_{k-1}(y)| + (M_F M_0 + 1)m_u/2 \\ &= \max_{x \in [0, x_k]} |\hat{v}'_{k-1}(x)| + (M_F M_0 + 1)m_u/2 \\ &\leq \max_{x \in [0, x_{k-1}]} |\hat{v}'_{k-1}(x)| + (M_F M_0 + 1)m_u/2 \\ &\leq (k-1)(M_F M_0 + 1)m_u/2 + (M_F M_0 + 1)m_u/2 = k(M_F M_0 + 1)m_u/2, \end{aligned} \quad (4.35)$$

just as we needed.

Step 2. As the second step, we quantitatively estimate the two auxiliary functions $u_k(m, x)$ and $\ell_k(m, x)$ when $k \geq K_0 + 1$ and $x \leq M_0 \epsilon_{k-1}(x)$. First we apply the integral representation (4.8) twice to get the identity

$$\int_{\epsilon_k(x)}^{\epsilon_{k-1}(x)} w f(w) dw = \frac{x}{k(k-1)} = \frac{1}{k-1} \int_0^{\epsilon_k(x)} w f(w) dw. \quad (4.36)$$

The monotonicity of the integrand $w f(w)$ from (4.17) implies a lower bound on the right integral in (4.36)

$$\int_{\epsilon_k(x)}^{\epsilon_{k-1}(x)} w f(w) dw \geq (\epsilon_{k-1}(x) - \epsilon_k(x)) \epsilon_k(x) f(\epsilon_k(x)),$$

as well as an upper bound on the left integral in (4.36)

$$\int_0^{\epsilon_k(x)} w f(w) dw \leq \epsilon_k^2(x) f(\epsilon_k(x)).$$

Taking these last two bounds back to (4.36) and simplify, we obtain the inequality

$$\epsilon_{k-1}(x) \leq \left(1 + \frac{1}{k-1}\right) \epsilon_k(x) \leq (1 + K_0^{-1}) \epsilon_k(x), \quad (4.37)$$

for all $k \geq K_0 + 1$.

Next, the integral representation of the consumption function (4.8) and the boundedness of the density f in (4.12) imply that

$$\frac{x}{k} = \int_0^{\epsilon_k(x)} w f(w) dw \geq m_\ell \int_0^{\epsilon_k(x)} w dw = \frac{1}{2} m_\ell \epsilon_k^2(x),$$

which in turn gives us an upper bound on the consumption function $\epsilon_k(x) \leq \sqrt{2x/(km_\ell)}$. This upper bound together with (4.37) yields that when $x \leq M_0 \epsilon_{k-1}(x)$, we have that

$$x \leq M_0 \epsilon_{k-1}(x) \leq M_0(1 + K_0^{-1}) \epsilon_k(x) \leq M_0(1 + K_0^{-1}) \sqrt{\frac{2x}{km_\ell}}.$$

Solving this inequality with respect to x , we obtain that

$$x \leq \frac{2M_0^2(1 - K_0^{-1})^2}{km_\ell},$$

and hence

$$x^{-3/4} \epsilon_k^{-1/4}(x) \geq x^{-3/4} \left(\sqrt{\frac{2x}{km_\ell}} \right)^{-1/4} \geq \frac{km_\ell}{2[M_0(1 + K_0^{-1})]^{7/4}}. \quad (4.38)$$

As a result, we have the lower bound

$$u_k(m, x) = \frac{1}{\epsilon_k(x)} + \frac{m}{x^{3/4} \epsilon_k^{1/4}(x)} \geq \frac{m}{x^{3/4} \epsilon_k^{1/4}(x)} \geq \frac{km m_\ell}{2[M_0(1 + K_0^{-1})]^{7/4}}.$$

Since $m \geq m_0$, in particular the second maximand in (4.31), we have that

$$u_k(m, x) \geq k(M_F M_0 + 1) m_u / 2,$$

which, together with (4.35), imply that $\hat{v}'_{k-1}(x) \leq u_k(m, x)$.

On the other hand, when $x \leq M_0 \epsilon_{k-1}(x)$, we have the upper bound

$$\ell_k(m, x) = \frac{1}{\epsilon_k(x)} - \frac{m}{x^{3/4} \epsilon_k^{1/4}(x)} = -\frac{m - (x/\epsilon_k(x))^{3/4}}{x^{3/4} \epsilon_k^{1/4}(x)} \leq -\frac{m - M_0^{3/4} [\epsilon_{k-1}(x)/\epsilon_k(x)]^{3/4}}{x^{3/4} \epsilon_k^{1/4}(x)}.$$

An application of (4.37) further implies that

$$\ell_k(m, x) \leq -\frac{m - [M_0(1 + K_0^{-1})]^{3/4}}{x^{3/4} \epsilon_k^{1/4}(x)}.$$

Therefore, when $m \geq m_0$, in particular the third maximand in (4.31), we have that

$$\ell_k(m, x) \leq -\frac{m_u(M_F M_0 + 1) [M_0(1 + K_0^{-1})]^{7/4}}{m_\ell x^{3/4} \epsilon_k^{1/4}(x)}.$$

Another application of (4.38) yields that

$$\ell_k(m, x) \leq -k(M_F M_0 + 1)m_u/2,$$

and therefore by (4.35) we know that $\ell_k(m, x) \leq \hat{v}_{k-1}'(x)$, and the proof is complete. \square

The second part of the induction argument for proving Lemma 4.1 is to estimate the derivative $\hat{v}_k'(x)$ in the other region $\Omega_k^c = [0, c] \setminus \Omega_k = \{x \in [0, c] : x > M_0 \epsilon_k(x)\}$. This estimation is given in the following proposition.

Proposition 4.8 (Intermediate step for induction). *Consider the two auxiliary functions*

$$u_k(m, x) = \frac{1}{\epsilon_k(x)} + \frac{m}{x^{3/4} \epsilon_k^{1/4}(x)}, \quad \text{and} \quad \ell_k(m, x) = \frac{1}{\epsilon_k(x)} - \frac{m}{x^{3/4} \epsilon_k^{1/4}(x)},$$

with $m \geq m_0$ where m_0 is defined in (4.31). If the item-weight distribution F belongs to the regular class, then for all $k \geq K_0$ and all $x \in \Omega_k^c \equiv \{x \in [0, c] : x > M_0 \epsilon_k(x)\}$, we have the following two inequalities

$$\begin{aligned} u_{k+1}(m, x) &\geq [1 - F(\epsilon_k(x))] u_k(m, x) + \int_0^{\epsilon_k(x)} u_k(m, x - w) f(w) dw \\ &\quad + \frac{1}{k \epsilon_k(x)} \left[1 - \int_0^{\epsilon_k(x)} u_k(m, x - w) dw \right], \end{aligned} \quad (4.39)$$

$$\begin{aligned} \ell_{k+1}(m, x) &\leq [1 - F(\epsilon_k(x))] \ell_k(m, x) + \int_0^{\epsilon_k(x)} \ell_k(m, x - w) f(w) dw \\ &\quad + \frac{1}{k \epsilon_k(x)} \left[1 - \int_0^{\epsilon_k(x)} \ell_k(m, x - w) dw \right]. \end{aligned} \quad (4.40)$$

Proof. With some simple rearrangements, inequality (4.39) is equivalent to

$$u_{k+1}(m, x) - u_k(m, x) + F(\epsilon_k(x)) u_k(m, x) - \int_0^{\epsilon_k(x)} u_k(m, x - w) \left\{ f(w) - \frac{1}{k \epsilon_k(x)} \right\} dw - \frac{1}{k \epsilon_k(x)} \geq 0.$$

Next, we plug in the representation $u_k(m, x) = \epsilon_k^{-1}(x) + mx^{-3/4}\epsilon_k^{-1/4}(x)$ and collect all the terms that result from $\epsilon_k^{-1}(x)$ as well as the last term $-(k\epsilon_k(x))^{-1}$ into one group $G_k(x)$, and all the rest terms into another group $H_k(x)$ to obtain that

$$\begin{aligned} & u_{k+1}(m, x) - u_k(m, x) + F(\epsilon_k(x))u_k(m, x) - \int_0^{\epsilon_k(x)} u_k(m, x-w) \left\{ f(w) - \frac{1}{k\epsilon_k(x)} \right\} dw - \frac{1}{k\epsilon_k(x)} \\ & = G_k(x) + mH_k(x), \end{aligned}$$

where the first group of terms are collected in

$$G_k(x) \equiv \frac{1}{\epsilon_{k+1}(x)} - \frac{1}{\epsilon_k(x)} - \frac{1}{k\epsilon_k(x)} \left\{ 1 - \int_0^{\epsilon_k(x)} \frac{1}{\epsilon_k(x-w)} dw \right\} - \int_0^{\epsilon_k(x)} \left\{ \frac{1}{\epsilon_k(x-w)} - \frac{1}{\epsilon_k(x)} \right\} f(w) dw,$$

and the second group of terms are collected in

$$\begin{aligned} H(k, x) \equiv & \frac{1}{x^{3/4}\epsilon_{k+1}^{1/4}(x)} - \frac{1}{x^{3/4}\epsilon_k^{1/4}(x)} + \frac{1}{k\epsilon_k(x)} \int_0^{\epsilon_k(x)} \frac{1}{(x-w)^{3/4}\epsilon_k^{1/4}(x-w)} dw \\ & - \int_0^{\epsilon_k(x)} \left\{ \frac{1}{(x-w)^{3/4}\epsilon_k^{1/4}(x-w)} - \frac{1}{x^{3/4}\epsilon_k^{1/4}(x)} \right\} f(w) dw. \end{aligned}$$

Therefore, the first inequality (4.39) is equivalent to $G_k(x) + mH_k(x) \geq 0$. Due to the symmetricity between the functions $\ell_k(m, x)$ and $u_k(m, x)$, similar algebra implies that the second inequality (4.40) is equivalent to $G_k(x) - mH_k(x) \leq 0$. Hence, the proof of the two inequalities (4.39) and (4.40) will complete once we can show that

$$m \cdot H_k(x) \geq |G_k(x)| \quad \text{for all } x \in \Omega_k^c.$$

The rest of the proof consists of three steps. The first step is to give an upper bound on $|G_k(x)|$; the second step is to give a positive lower bound on $H_k(x)$; and the last step is to show that the choice of $m \geq m_0$ guarantees that the lower bound on $mH_k(x)$ always dominates the upper bound on $|G_k(x)|$.

Step 1. We will show that $|G_k(x)| \leq (M_F + 1)(kx)^{-1}$.

Recall that the function $G_k(x)$ is defined as

$$\begin{aligned} G_k(x) &= \left\{ \frac{1}{\epsilon_{k+1}(x)} - \frac{1}{\epsilon_k(x)} \right\} - \frac{1}{k\epsilon_k(x)} \left\{ 1 - \int_0^{\epsilon_k(x)} \frac{1}{\epsilon_k(x-w)} dw \right\} \\ &\quad - \int_0^{\epsilon_k(x)} \left\{ \frac{1}{\epsilon_k(x-w)} - \frac{1}{\epsilon_k(x)} \right\} f(w) dw \\ &= I_G + II_G + III_G. \end{aligned}$$

Towards the upper bound on $|G_k(x)|$, we start with relaxing the first piece through the integral representation

$$I_G = \frac{1}{\epsilon_{k+1}(x)} - \frac{1}{\epsilon_k(x)} = \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} \frac{dw}{w^2} = \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} \frac{wf(w)}{w^3 f(w)} dw.$$

The monotonicity of $wf(w)$ from (4.17) tells us that $w^3 f(w)$ is also monotone increasing, and hence in the right-most integral of this last equation, if we replace the denominator $w^3 f(w)$ of the integrand with its upper bound $\epsilon_k^3(x) f(\epsilon_k(x))$, we get a lower bound

$$I_G \geq \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} \frac{wf(w)}{\epsilon_k^3(x) f(\epsilon_k(x))} dw = \frac{x}{k(k+1)\epsilon_k^3(x) f(\epsilon_k(x))}, \quad (4.41)$$

where the last identity is due to the representation of the consumption function (4.8). On the other hand, we have another integral representation of the first piece I_G given by

$$I_G = \frac{1}{\epsilon_{k+1}(x)} - \frac{1}{\epsilon_k(x)} = \frac{\epsilon_k(x) - \epsilon_{k+1}(x)}{\epsilon_k(x)\epsilon_{k+1}(x)} = \frac{1}{\epsilon_k(x)\epsilon_{k+1}(x)} \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} 1 dw.$$

In this last right-hand side, if we use the monotonicity of $wf(w)$ from (4.17) once again to replace the integrand 1 with its upper bound $wf(w)/(\epsilon_{k+1}(x)f(\epsilon_{k+1}(x)))$, we find

$$I_G \leq \frac{1}{\epsilon_k(x)\epsilon_{k+1}(x)} \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} \frac{wf(w)}{\epsilon_{k+1}(x)f(\epsilon_{k+1}(x))} dw$$

Then we apply in turn the consumption function representation (4.8) and the monotone relation (4.24) to get

$$I_G \leq \frac{x}{k(k+1)\epsilon_k(x)\epsilon_{k+1}^2(x)f(\epsilon_{k+1}(x))} \leq \frac{(k+1)x}{k^3\epsilon_k^3(x)f(\epsilon_k(x))}. \quad (4.42)$$

Putting both the lower bound (4.41) and the upper bound (4.42) together, we have the two-sided bound

$$\frac{x}{k(k+1)\epsilon_k^3(x)f(\epsilon_k(x))} \leq I_G \leq \frac{(k+1)x}{k^3\epsilon_k^3(x)f(\epsilon_k(x))}. \quad (4.43)$$

For the second piece II_G of the function $G_k(x)$, we know from (4.10) the closed-form integral representation

$$\int \frac{1}{\epsilon_k(w)} dw = kF(\epsilon_k(w)),$$

which in turn gives us the closed-form identity

$$II_G = -\frac{1}{k\epsilon_k(x)} \left\{ 1 - \int_0^{\epsilon_k(x)} \frac{1}{\epsilon_k(x-w)} dw \right\} = \frac{1}{k\epsilon_k(x)} \{kF(\epsilon_k(x)) - kF(\epsilon_k(x - \epsilon_k(x))) - 1\}.$$

Now, if we apply (4.23) with $w = \epsilon_k(x)$, we get the inequality

$$0 \leq kF(\epsilon_k(x)) - kF(\epsilon_k(x - \epsilon_k(x))) - 1 \leq \frac{k\epsilon_k^2(x)F(\epsilon_k(x))}{x^2},$$

which further gives us the two-sided inequality on the second piece

$$0 \leq II_G \leq \frac{\epsilon_k(x)F(\epsilon_k(x))}{x^2}. \quad (4.44)$$

For the third piece

$$III_G = -\int_0^{\epsilon_k(x)} \left\{ \frac{1}{\epsilon_k(x-w)} - \frac{1}{\epsilon_k(x)} \right\} f(w) dw,$$

since $x \in \Omega_k^c$, we know that $\epsilon_k(x) \leq M_0^{-1}x \leq x/2$ as $M_0 \geq 2$, so (4.25) implies that for all $w \in [0, \epsilon_k(x)]$ the first factor of the integrand of III_G satisfies the two-sided bound

$$\frac{w}{k\epsilon^3(x)f(\epsilon_k(x))} \leq \frac{1}{\epsilon_k(x-w)} - \frac{1}{\epsilon_k(x)} \leq \frac{w}{k\epsilon^3(x)f(\epsilon_k(x))} + M_F \cdot \frac{w^2}{x^2\epsilon_k(x)}.$$

Then we multiply this last equation with $f(w)$, integrate over $w \in [0, \epsilon_k(x)]$ and use the representation (4.8) to get

$$\begin{aligned} \frac{x}{k^2\epsilon^3(x)f(\epsilon_k(x))} &\leq \int_0^{\epsilon_k(x)} \left\{ \frac{1}{\epsilon_k(x-w)} - \frac{1}{\epsilon_k(x)} \right\} f(w) dw \\ &\leq \frac{x}{k^2\epsilon^3(x)f(\epsilon_k(x))} + \frac{M_F}{x^2\epsilon_k(x)} \int_0^{\epsilon_k(x)} w^2 f(w) dw. \end{aligned}$$

In the rightmost integral of this last inequality, if we replace the integrand $w^2 f(w)$ with its upper bound $\epsilon_k(x)w f(w)$ and use (4.8) once again, we obtain that

$$\frac{x}{k^2\epsilon^3(x)f(\epsilon_k(x))} \leq \int_0^{\epsilon_k(x)} \left\{ \frac{1}{\epsilon_k(x-w)} - \frac{1}{\epsilon_k(x)} \right\} f(w) dw \leq \frac{x}{k^2\epsilon^3(x)f(\epsilon_k(x))} + \frac{M_F}{kx}.$$

By multiplying this last equation with (-1) to reverse the direction, we obtain that

$$-\frac{x}{k^2\epsilon^3(x)f(\epsilon_k(x))} - \frac{M_F}{kx} \leq III_G \leq -\frac{x}{k^2\epsilon^3(x)f(\epsilon_k(x))}. \quad (4.45)$$

Lastly, we combine all the three inequalities (4.43), (4.44) and (4.45) to get the final estimation

$$-\frac{x}{k^2(k+1)\epsilon_k^3(x)f(\epsilon_k(x))} - \frac{M_F}{kx} \leq G_k(x) \leq \frac{x}{k^3\epsilon_k^3(x)f(\epsilon_k(x))} + \frac{M_F}{kx}.$$

Taking absolute value on this last equation, we obtain that

$$|G_k(x)| \leq \frac{x}{k^3\epsilon_k^3(x)f(\epsilon_k(x))} + \frac{M_F}{kx}. \quad (4.46)$$

If we rewrite the first term in this last right-hand side as

$$\frac{x}{k^3\epsilon_k^3(x)f(\epsilon_k(x))} = \frac{1}{kx} \left(\frac{x}{k\epsilon_k^2(x)f(\epsilon_k(x))} \right) \frac{x}{k\epsilon_k(x)},$$

and apply (4.22) to bound the second factor from above by 1, and apply (4.21) to bound the third factor from above by $F(\epsilon_k(x))$, we get

$$\frac{x}{k^3\epsilon_k^3(x)f(\epsilon_k(x))} \leq \frac{F(\epsilon_k(x))}{kx} \leq \frac{1}{kx}.$$

Putting this back to (4.46), we get the final estimation

$$|G_k(x)| \leq \frac{M_F + 1}{kx}. \quad (4.47)$$

Step 2. We will show that $H_k(x) \geq \{8kx^{3/4}\epsilon_k^{1/4}(x)\}^{-1}$.

Recall that the function $H_k(x)$ is defined as

$$\begin{aligned} H_k(x) &= \frac{1}{x^{3/4}\epsilon_{k+1}^{1/4}(x)} - \frac{1}{x^{3/4}\epsilon_k^{1/4}(x)} + \frac{1}{k\epsilon_k(x)} \int_0^{\epsilon_k(x)} \frac{1}{(x-w)^{3/4}\epsilon_k^{1/4}(x-w)} dw \\ &\quad - \int_0^{\epsilon_k(x)} \left\{ \frac{1}{(x-w)^{3/4}\epsilon_k^{1/4}(x-w)} - \frac{1}{x^{3/4}\epsilon_k^{1/4}(x)} \right\} f(w) dw \\ &= I_H + II_H + III_H. \end{aligned}$$

Similar to what we did to the function $G_k(x)$, we exploit the monotonicity of $wf(w)$ to estimation the first piece I_H . In fact, we have the integral representation

$$I_H = x^{-3/4}\epsilon_{k+1}^{-1/4}(x) - x^{-3/4}\epsilon_k^{-1/4}(x) = \frac{1}{4x^{3/4}} \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} w^{-5/4} dw.$$

The monotonicity of $wf(w)$ from (4.17) implies that $w^{9/4}f(w)$ is also increasing, so for all $w \in [\epsilon_{k+1}(x), \epsilon_k(x)]$,

$$w^{-5/4} \geq \frac{wf(w)}{\epsilon_k^{9/4}(x)f(\epsilon_k(x))},$$

and as a result,

$$I_H \geq \frac{1}{4x^{3/4}} \int_{\epsilon_{k+1}(x)}^{\epsilon_k(x)} \frac{wf(w)}{\epsilon_k^{9/4}(x)f(\epsilon_k(x))} dw = \frac{x^{1/4}}{4k(k+1)\epsilon_k^{9/4}(x)f(\epsilon_k(x))}, \quad (4.48)$$

where the last equality is due to (4.8).

The second piece

$$II_H = \frac{1}{k\epsilon_k(x)} \int_0^{\epsilon_k(x)} \frac{1}{(x-w)^{3/4}\epsilon_k^{1/4}(x-w)} dw$$

has a crude lower bound if we replace the integrand $((x-w)^{3/4}\epsilon_k^{1/4}(x-w))^{-1}$ with its lower bound $(x^{3/4}\epsilon_k^{1/4}(x))^{-1}$ for all $w \in [0, \epsilon_k(x)]$. That is,

$$II_H \geq \frac{1}{k\epsilon_k(x)} \int_0^{\epsilon_k(x)} \frac{1}{x^{3/4}\epsilon_k^{1/4}(x)} dw = \frac{1}{kx^{3/4}\epsilon_k^{1/4}(x)}. \quad (4.49)$$

For the third piece

$$III_H = - \int_0^{\epsilon_k(x)} \left((x-w)^{-3/4}\epsilon_k^{-1/4}(x-w) - x^{-3/4}\epsilon_k^{-1/4}(x) \right) f(w) dw,$$

we take $\alpha = 1/4$ in (4.26) to get the bound on the first factor of the integrand

$$(x-w)^{-3/4}\epsilon_k^{-1/4}(x-w) - x^{-3/4}\epsilon_k^{-1/4}(x) \leq \left(\frac{3}{4x^{7/4}\epsilon_k^{1/4}(x)} + \frac{1}{4kx^{3/4}\epsilon_k^{9/4}(x)f(\epsilon_k(x))} \right) w + \frac{M_F w^2}{x^{11/4}\epsilon_k^{1/4}(x)}.$$

As a result, if we multiply this last inequality with $(-f(w))$ and integrate over $w \in [0, \epsilon_k(x)]$, we obtain that

$$III_H \geq - \int_0^{\epsilon_k(x)} \left(\frac{3}{4x^{7/4}\epsilon_k^{1/4}(x)} + \frac{1}{4kx^{3/4}\epsilon_k^{9/4}(x)f(\epsilon_k(x))} \right) wf(w) + \frac{M_F w^2 f(w)}{x^{11/4}\epsilon_k^{1/4}(x)} dw.$$

In the last summand of the integrand of this last right-hand side, if we replace $w^2 f(w)$ with its upper bound $\epsilon_k(x)wf(w)$, we obtain a further lower bound

$$III_H \geq - \int_0^{\epsilon_k(x)} \left(\frac{3}{4x^{7/4}\epsilon_k^{1/4}(x)} + \frac{1}{4kx^{3/4}\epsilon_k^{9/4}(x)f(\epsilon_k(x))} + \frac{M_F \epsilon_k^{3/4}(x)}{x^{11/4}} \right) wf(w) dw.$$

Another application of (4.8) yields that

$$III_H \geq -\frac{3}{4kx^{3/4}\epsilon_k^{1/4}(x)} - \frac{x^{1/4}}{4k^2\epsilon_k^{9/4}(x)f(\epsilon_k(x))} - \frac{M_F\epsilon_k(x)}{kx^{7/4}}. \quad (4.50)$$

Combining all three inequalities (4.48), (4.49) and (4.50) together, we obtain the overall estimation

$$H_k(x) \geq \left(\frac{k}{4(k+1)} - \frac{1}{4}\right) \frac{x^{1/4}}{k^2\epsilon_k^{9/4}(x)f(\epsilon_k(x))} + \left(\frac{1}{4} - \frac{M_F\epsilon_k(x)}{x}\right) \frac{1}{kx^{3/4}\epsilon_k^{1/4}(x)}.$$

From the left inequality in (4.22), we know that

$$\frac{x^{1/4}}{k^2\epsilon_k^{9/4}(x)f(\epsilon_k(x))} \leq \frac{1}{kx^{3/4}\epsilon_k^{1/4}(x)},$$

and hence

$$H_k(x) \geq \left(\frac{k}{4(k+1)} - \frac{M_F\epsilon_k(x)}{x}\right) \frac{1}{kx^{3/4}\epsilon_k^{1/4}(x)}.$$

The domain $\Omega_k^c = \{x \in [0, c] : x > M_0\epsilon_k(x)\}$ of x yields the further lower bound

$$H_k(x) \geq \left(\frac{k}{4(k+1)} - \frac{M_F}{M_0}\right) \frac{1}{kx^{3/4}\epsilon_k^{1/4}(x)}.$$

The choice of M_0 from (4.30) and K_0 from (4.20) guarantees us that the first factor in this last right-hand side is at least $\frac{k}{4(k+1)} - \frac{M_F}{M_0} \geq 1/8$, and so

$$H(k, x) \geq \frac{1}{8kx^{3/4}\epsilon_k^{1/4}(x)}. \quad (4.51)$$

Step 3. We will show that $m \geq m_0$ implies that $|G_k(x)| \leq m \cdot H_k(x)$.

In fact, from the choice of the constant m_0 , specifically the first maximand in (4.31), we know that when $m \geq m_0$,

$$\left(\frac{8(M_F + 1)}{m}\right)^4 \leq \left(\frac{8(M_F + 1)}{m_0}\right)^4 \leq M_0.$$

Since $x \in \Omega_k^c = \{x \in [0, c] : x > M_0\epsilon_k(x)\}$, we further have that

$$\frac{8(M_F + 1)}{m} \leq \frac{x^{1/4}}{\epsilon_k^{1/4}(x)},$$

which implies that

$$\frac{M_F + 1}{kx} \leq \frac{m}{8kx^{3/4}\epsilon_k^{1/4}(x)}.$$

This last inequality together with the two estimations (4.47) and (4.51) imply that

$$|G_k(x)| \leq m \cdot H_k(x),$$

just as we needed. The proof is complete. \square

With the last two propositions, we are now ready to prove the key lemma of this section.

Proof of Lemma 4.1. With the constant K_0 defined in (4.20), we consider an auxiliary function

$$\delta(x) \equiv \left| \widehat{v}'_{K_0-1}(x) \cdot x^{3/4}\epsilon_{K_0}^{1/4}(x) - \left(\frac{x}{\epsilon_{K_0}(x)} \right)^{3/4} \right|, \quad x \in (0, c].$$

The limit $\lim_{x \rightarrow 0} \epsilon_{K_0}(x) = 0$ together with the boundedness of density f from (4.12) enable us to apply L'Hospital's rule with (4.10) to get the limit

$$\lim_{x \rightarrow 0^+} \frac{x}{\epsilon_{K_0}(x)} = \lim_{x \rightarrow 0^+} K_0 \epsilon_{K_0}(x) f(\epsilon_{K_0}(x)) = 0.$$

As a result, the function $\delta(x)$ defined on $(0, c]$ can be extended to the entire interval $[0, c]$, hence this continuous function $\delta(x)$ attains its finite maximum over the compact domain $[0, c]$. Then we can define constant m as

$$m \equiv \max \left\{ \max_{x \in [0, c]} \delta(x), m_0 \right\} < \infty. \quad (4.52)$$

This last definition of m gives us the induction base, that is,

$$\ell_{K_0}(m, x) \leq \widehat{v}'_{K_0-1}(x) \leq u_{K_0}(m, x), \quad \text{for all } x \in [0, c].$$

Now we proceed to the induction step. Assume that

$$\ell_k(m, x) \leq \widehat{v}'_{k-1}(x) \leq u_k(m, x), \quad \text{for all } x \in [0, c]. \quad (4.53)$$

Recall the k -th partition of the entire domain $[0, c] = \Omega_k \cup \Omega_k^c$, where $\Omega_k = \{x \in [0, c] : x \leq M_0 \epsilon_k(x)\}$. From Proposition 4.7 we know that for all $x \in \Omega_k$,

$$\ell_{k+1}(m, x) \leq \widehat{v}'_k(x) \leq u_{k+1}(m, x).$$

Then to complete the proof, it suffices to show that (4.53) implies that for all $x \in \Omega_k^c$,

$$\ell_{k+1}(m, x) \leq \widehat{v}'_k(x) \leq u_{k+1}(m, x).$$

The integral representation (4.8) together with the boundedness of the density f from (4.12) yield the chain relation

$$\frac{x}{k} = \int_0^{\epsilon_k(x)} w f(w) dw \leq m_u \int_0^{\epsilon_k(x)} w dw = \frac{m_u \epsilon_k^2(x)}{2},$$

which in turn implies that

$$\frac{1}{k \epsilon_k(x)} \leq \frac{m_u \epsilon_k(x)}{2x}.$$

When $x \in \Omega_k^c = \{x \in [0, c] : x > M_0 \epsilon_k(x)\}$, this last inequality gives us a further upper bound

$$\frac{1}{k \epsilon_k(x)} \leq \frac{m_u}{2M_0}.$$

This last upper bound together with the boundedness of the density f from (4.12) as well as the definition of the constant M_0 in (4.30) imply that for all $w \in [0, \epsilon_k(x)]$,

$$f(w) \geq m_\ell \geq \frac{m_u}{2M_0} \geq \frac{1}{k \epsilon_k(x)}. \quad (4.54)$$

Now, let us turn back to the recursion of the derivative $\widehat{v}'_k(x)$. For all $x \in \Omega_k^c$, we know that $x > M_0 \epsilon_k(x) \geq \epsilon_k(x)$, so the recursion of $\widehat{v}'_k(x)$ is given by

$$\begin{aligned} \widehat{v}'_k(x) &= (1 - F(\epsilon_k(x))) v'_{k-1}(x) + \int_0^{\epsilon_k(x)} v'_{k-1}(x - w) f(w) dw + \frac{1}{k \epsilon_k(x)} \left(1 - \int_0^{\epsilon_k(x)} v'_{k-1}(x - w) dw \right) \\ &= (1 - F(\epsilon_k(x))) v'_{k-1}(x) + \int_0^{\epsilon_k(x)} v'_{k-1}(x - w) \left(f(w) - \frac{1}{k \epsilon_k(x)} \right) dw + \frac{1}{k \epsilon_k(x)}. \end{aligned}$$

Then by the right inequality of our induction assumption (4.53) and the relation (4.54), we know that

$$\begin{aligned} \widehat{v}'_k(x) &\leq (1 - F(\epsilon_k(x))) u_k(m, x) + \int_0^{\epsilon_k(x)} u_k(m, x - w) \left(f(w) - \frac{1}{k \epsilon_k(x)} \right) dw + \frac{1}{k \epsilon_k(x)} \\ &= (1 - F(\epsilon_k(x))) u_k(m, x) + \int_0^{\epsilon_k(x)} u_k(m, x - w) f(w) dw + \frac{1}{k \epsilon_k(x)} \left[1 - \int_0^{\epsilon_k(x)} u_k(m, x - w) dw \right], \end{aligned}$$

so we can apply (4.39) to this last right-hand side to obtain that

$$\widehat{v}'_k(x) \leq u_{k+1}(m, x).$$

Parallely, by the left inequality of our induction assumption (4.53) and the relation (4.54), we have the lower bound

$$\begin{aligned} \widehat{v}'_k(x) &\geq (1 - F(\epsilon_k(x))) \ell_k(m, x) + \int_0^{\epsilon_k(x)} \ell_k(m, x - w) \left(f(w) - \frac{1}{k\epsilon_k(x)} \right) dw + \frac{1}{k\epsilon_k(x)} \\ &= (1 - F(\epsilon_k(x))) \ell_k(m, x) + \int_0^{\epsilon_k(x)} \ell_k(m, x - w) f(w) dw + \frac{1}{k\epsilon_k(x)} \left[1 - \int_0^{\epsilon_k(x)} \ell_k(m, x - w) dw \right]. \end{aligned}$$

Then we apply (4.40) to this last right-hand side to obtain that

$$\widehat{v}'_k(x) \geq \ell_{k+1}(m, x).$$

To sum up, the induction assumption (4.53) does imply the two-sided inequality

$$\ell_{k+1}(m, x) \leq \widehat{v}'_k(x) \leq u_{k+1}(m, x), \quad \text{for all } x \in \Omega_k^c,$$

and therefore the inequality

$$\ell_k(m, x) \leq \widehat{v}'_{k-1}(x) \leq u_k(m, x) \quad \text{for all } x \in [0, c]$$

hold for all $k \geq K_0$.

We conclude this section with a corollary of Lemma 4.1 which states that the difference in the value function $|\widehat{v}_{k-1}(x) - \widehat{v}_{k-1}(x - w)|$ within the interval $w \in [0, \widehat{h}_k(x)]$ is uniformly bounded by constant that is independent of k .

Corollary 4.1 (Bounded difference in value functions within selection interval). *There exists a constant $M \equiv \max\{4M_F, (m+1)(M_F/2+1)\}$ that depends only on F , such that for all $k \geq 1$ and all $x \in [0, c]$,*

$$|\widehat{v}_{k-1}(x) - \widehat{v}_{k-1}(x - w)| \leq M \quad \text{for all } w \in [0, \min\{x, \epsilon_k(x)\}].$$

Proof. Take any $k \geq 1$ and $x \in [0, c]$ fixed. There are only two possible scenarios: (i) $x \leq 2\epsilon_k(x)$ or (ii) $x > 2\epsilon_k(x)$.

Case (i). If $x \leq 2\epsilon_k(x)$.

We apply Theorem 2.1 to obtain that, as a feasible online policy, $\hat{\pi}_n(c)$, its value function has upper bound for all $y \in [0, x]$,

$$\hat{v}_{k-1}(y) \leq (k-1)F(\epsilon_{k-1}(y)) \leq (k-1)F(\epsilon_{k-1}(x)) = \frac{(k-1)\epsilon_{k-1}(x)F(\epsilon_{k-1}(x))}{x} \frac{x}{\epsilon_{k-1}(x)}.$$

If we apply (4.21) to the first factor of this last right-hand side, and note that $x \leq 2\epsilon_k(x) \leq 2\epsilon_{k-1}(x)$, we find that

$$\hat{v}_{k-1}(y) \leq 2M_F \quad \text{for all } y \in [0, x].$$

As a result, we obtain that if $x \leq 2\epsilon_k(x)$, then for all $w \in [0, \min\{x, \epsilon_k(x)\}] \subset [0, x]$,

$$|\hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w)| \leq 4M_F.$$

Case (ii). If $x > 2\epsilon_k(x)$.

In this case, for all $w \in [0, \epsilon_k(x)]$, we have the relation $x-w \geq x-\epsilon_k(x) > \epsilon_k(x) \geq \epsilon_k(x-w)$, so

$$\frac{1}{(x-w)^{3/4}\epsilon_k^{1/4}(x-w)} = \left(\frac{\epsilon_k(x-w)}{x-w}\right)^{3/4} \frac{1}{\epsilon_k(x-w)} \leq \frac{1}{\epsilon_k(x-w)}.$$

This last inequality together with Lemma 4.1 imply that for all $y \in [x-\epsilon_k(x), x]$,

$$|\hat{v}'_{k-1}(y)| \leq \frac{1}{\epsilon_k(y)} + \frac{m}{y^{3/4}\epsilon_k^{1/4}(y)} \leq \frac{m+1}{\epsilon_k(y)}.$$

Therefore, by the fundamental theorem of calculus, we have that for all $w \in [0, \epsilon_k(x)]$,

$$\begin{aligned} & |\hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w)| \\ &= \left| \int_{x-w}^x \hat{v}'_{k-1}(y) dy \right| \leq \int_{x-w}^x |\hat{v}'_{k-1}(y)| dy \leq \int_{x-\epsilon_k(x)}^x |\hat{v}'_{k-1}(y)| dy \\ &\leq \int_{x-\epsilon_k(x)}^x \frac{m+1}{\epsilon_k(y)} dy = (m+1) [kF(\epsilon_k(x)) - kF(\epsilon_k(x-\epsilon_k(x)))], \end{aligned}$$

where in the last identity we have used $\int \epsilon_k^{-1}(y) dy = kF(\epsilon_k(y))$ due to (4.10). If we set the parameter w in (4.23) to $w = \epsilon_k(x)$, we obtain that

$$kF(\epsilon_k(x)) - kF(\epsilon_k(x-\epsilon_k(x))) \leq \frac{kF(\epsilon_k(x))\epsilon_k^2(x)}{x^2} + 1 = \frac{kF(\epsilon_k(x))\epsilon_k(x)}{x} \frac{\epsilon_k(x)}{x} + 1.$$

Since $x > 2\epsilon_k(x)$, if we apply (4.21) once again to bound the first factor in this last right-hand side, we find that

$$kF(\epsilon_k(x)) - kF(\epsilon_k(x-\epsilon_k(x))) \leq M_F/2 + 1.$$

Therefore, we obtain the final bound that for all $w \in [0, \epsilon_k(x)]$,

$$|\hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w)| \leq (m+1)(M_F/2 + 1).$$

To sum up, in either case, the difference in the value function $|\hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w)|$ is bounded by constant that depends only on F , so the proof is complete. \square

4.3 Variance Asymptotics with Remainder Bounds

Establishing the asymptotics of $\text{Var} \left\{ \hat{N}_n(c) \right\}$ as $n \rightarrow \infty$ is one of the most important steps towards proving the weak convergence in Theorem 4.1. Due to the inter-dependency generated by the sequential decisions, we use martingales to analyze such inter-dependency.

4.3.1 Doob's Martingale and Conditional Variance of Martingale Differences

With the natural filtration $\{\mathcal{F}_t : t = 0, 1, \dots, n\}$, we define Doob's martingale as

$$M_t \equiv \mathbb{E} \left[\hat{N}_n(c) \mid \mathcal{F}_{t-1} \right], \quad \text{for } t \in [n].$$

Then we also have the associated martingale differences

$$d_t = M_t - M_{t-1}, \quad \text{for all } t \in [n].$$

Towards proving the martingale central limit theorem, it is an important intermediate step to tightly estimate the variance of the martingale differences. To this end, we consider the \mathcal{F}_t -measurable conditional variance

$$\omega_{n-t}(X_t) \equiv \text{Var} \left\{ \sum_{j=t+1}^n d_j \mid \mathcal{F}_t \right\}. \quad (4.55)$$

Then by the orthogonality of martingale differences, we also have the equivalent representation

$$\omega_{n-t}(X_t) = \sum_{j=t+1}^n \mathbb{E} [d_j^2 \mid \mathcal{F}_t].$$

If we denote G_t as the event of selecting the t -th item under policy $\hat{\pi}$, that is,

$$G_t \equiv \left\{ \omega : W_t(\omega) \in [0, \hat{h}_{n-t+1}(X_{t-1}(\omega))] \right\},$$

then for each $t \in [n]$, we have the decomposition of the martingale difference $d_t = A_t + B_t$, where the two quantities are given by

$$\begin{aligned} A_t &\equiv \hat{v}_{n-t}(X_{t-1}) - \hat{v}_{n-t+1}(X_{t-1}), \\ B_t &\equiv \{1 + \hat{v}_{n-t}(X_{t-1} - W_t) - \hat{v}_{n-t}(X_{t-1})\} \mathbb{1}\{G_t\}. \end{aligned}$$

By Corollary 4.1 we know that $\max\{\|A_t\|_\infty, \|B_t\|_\infty\} \leq M + 1$, so we have the uniform bound

$$\|d_t\|_\infty = \|A_t + B_t\|_\infty \leq 2(M + 1) \quad \text{for all } t \in [n]. \quad (4.56)$$

The zero-mean property of martingale difference $\mathbb{E}[d_t \mid \mathcal{F}_{t-1}] = 0$ translates into $A_t + \mathbb{E}[B_t \mid \mathcal{F}_{t-1}] = 0$, and hence

$$\mathbb{E}[d_t^2 \mid \mathcal{F}_{t-1}] = \mathbb{E}[B_t^2 \mid \mathcal{F}_{t-1}] - A_t^2. \quad (4.57)$$

The two quantities in this last right-hand side both have integral representations. From the value function recursion (4.6), we know that

$$A_t = - \int_0^{\hat{h}_{n-t+1}(X_{t-1})} \{1 + \hat{v}_{n-t}(X_{t-1} - w) - \hat{v}_{n-t}(X_{t-1})\} f(w) dw. \quad (4.58)$$

On the other hand, we have the integral representation of the conditional expectation

$$\mathbb{E}[B_t^2 \mid \mathcal{F}_{t-1}] = \int_0^{\hat{h}_{n-t+1}(X_{t-1})} \{1 + \hat{v}_{n-t}(X_{t-1} - w) - \hat{v}_{n-t}(X_{t-1})\}^2 f(w) dw. \quad (4.59)$$

Summing up identity (4.57) from $j = t + 1$ to n and taking one more conditional expectation with respect to \mathcal{F}_t , we find that

$$\omega_{n-t}(X_t) = \sum_{j=t+1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_t] = \sum_{j=t+1}^n \mathbb{E}[B_j^2 \mid \mathcal{F}_t] - \sum_{j=t+1}^n \mathbb{E}[A_j^2 \mid \mathcal{F}_t]. \quad (4.60)$$

As we will see shortly, the summation $\sum_{j=t+1}^n \mathbb{E}[A_j^2 \mid \mathcal{F}_t]$ in this last identity is negligible compared to the other summation $\sum_{j=t+1}^n \mathbb{E}[B_j^2 \mid \mathcal{F}_t]$, which drives most of the fluctuation of the selection process. Fortunately, with the help of Lemma 4.1, we are able to tightly approximate $\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}]$ by replacing the integrand of the right-hand side of (4.59),

$$\{1 + \hat{v}_{n-t}(X_{t-1} - w) - \hat{v}_{n-t}(X_{t-1})\}^2 \quad \text{for } w \in [0, \hat{h}_{n-t+1}(X_{t-1})],$$

with the square of the linear counterpart

$$\left(1 - \frac{w}{\widehat{h}_{n-t+1}(X_{t-1})}\right)^2 \quad \text{for } w \in [0, \widehat{h}_{n-t+1}(X_{t-1})].$$

Such approximation is then used to estimate the conditional variance $\omega_{n-t}(\widehat{X}_t)$.

4.3.2 Conditional Variance Bounds

The conditional variance $\omega_{n-t}(X_t)$ is estimated in the following lemma with upper bounds on the remainders.

Lemma 4.2 (Conditional variance bounds). *There exists constant M that depends only on F and c such that for all $t \in [n]$,*

$$|\omega_{n-t}(X_t) - \gamma(r)v_{n-t}(X_t)| \leq M \left\{ [(n-t)F(\epsilon_{n-t}(X_t))]^{1/4} + \log(n-t) + 1 \right\}. \quad (4.61)$$

The proof of this last lemma relies on an intermediate step that approximates the important quantity

$$\int_0^{\widehat{h}_k(x)} \{1 + \widehat{v}_{k-1}(x-w) - \widehat{v}_{k-1}(x)\}^2 f(w) dw$$

by replacing the integrand with the corresponding linear function. The error of this approximation is well controlled with the help of Lemma 4.1, as stated in the following proposition.

Proposition 4.9 (Linear approximation of integrand). *There exists a constant M that depends only on F , such that for all $k \geq K_0$ with K_0 defined as in (4.20), we have that*

$$\begin{aligned} & \left| \int_0^{\widehat{h}_k(x)} \{1 + \widehat{v}_{k-1}(x-w) - \widehat{v}_{k-1}(x)\}^2 f(w) dw - \int_0^{\widehat{h}_k(x)} \left(1 - \frac{w}{\widehat{h}_k(x)}\right)^2 f(w) dw \right| \\ & \leq M \left(\frac{F^{1/4}(\epsilon_k(x))}{k^{3/4}} + \frac{1}{k} \right). \end{aligned} \quad (4.62)$$

Proof. The left-hand side of (4.62) has the simple upper bound

$$\begin{aligned} & \left| \int_0^{\widehat{h}_k(x)} \{1 + \widehat{v}_{k-1}(x-w) - \widehat{v}_{k-1}(x)\}^2 f(w) dw - \int_0^{\widehat{h}_k(x)} \left(1 - \frac{w}{\widehat{h}_k(x)}\right)^2 f(w) dw \right| \\ & \leq \sup_{w \in [0, \widehat{h}_k(x)]} \left| 2 + \widehat{v}_{k-1}(x-w) - \widehat{v}_{k-1}(x) - \frac{w}{\widehat{h}_k(x)} \right| \int_0^{\widehat{h}_k(x)} \left| \widehat{v}_{k-1}(x) - \widehat{v}_{k-1}(x-w) - \frac{w}{\widehat{h}_k(x)} \right| f(w) dw. \end{aligned}$$

From Corollary 4.1 we have the uniform bound for all $x \in [0, c]$ and all $k \geq 1$,

$$\sup_{w \in [0, \hat{h}_k(x)]} |\hat{v}_{k-1}(x-w) - \hat{v}_{k-1}(x)| \leq M. \quad (4.63)$$

On the other hand, we also have the simple bound on the linear function

$$1 \leq 2 - \frac{w}{\hat{h}_k(x)} \leq 2, \quad \text{for all } w \in [0, \hat{h}_k(x)].$$

Thus,

$$\sup_{w \in [0, \hat{h}_k(x)]} \left| 2 + \hat{v}_{k-1}(x-w) - \hat{v}_{k-1}(x) - \frac{w}{\hat{h}_k(x)} \right| \leq M + 2,$$

and as a result, the left-hand side of (4.62) is bounded by

$$\begin{aligned} & \left| \int_0^{\hat{h}_k(x)} \{1 + \hat{v}_{k-1}(x-w) - \hat{v}_{k-1}(x)\}^2 f(w) dw - \int_0^{\hat{h}_k(x)} \left(1 - \frac{w}{\hat{h}_k(x)}\right)^2 f(w) dw \right| \\ & \leq (M + 2) \int_0^{\hat{h}_k(x)} \left| \hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w) - \frac{w}{\hat{h}_k(x)} \right| f(w) dw. \end{aligned}$$

Next, we recall the definition of the threshold function $\hat{h}_k(x) = \min\{x, \epsilon_k(x)\}$ from (4.4) and proceed the analysis for the only two possible scenarios: (i) $\hat{h}_k(x) = x$ and (ii) $\hat{h}_k(x) = \epsilon_k(x)$ separately.

Case (i) $\hat{h}_k(x) = x \leq \epsilon_k(x)$. The uniform bound (4.63) implies that the integrand of this last right-hand side satisfies that

$$\begin{aligned} & \sup_{w \in [0, \hat{h}_k(x)]} \left| \hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w) - \frac{w}{\hat{h}_k(x)} \right| \\ & \leq \sup_{w \in [0, \hat{h}_k(x)]} |\hat{v}_{k-1}(x-w) - \hat{v}_{k-1}(x)| + \sup_{w \in [0, \hat{h}_k(x)]} \left| \frac{w}{\hat{h}_k(x)} \right| \leq M + 1. \end{aligned}$$

As a result, the left-hand side of (4.62) is further bounded from above by

$$\begin{aligned} & \left| \int_0^{\hat{h}_k(x)} \{1 + \hat{v}_{k-1}(x-w) - \hat{v}_{k-1}(x)\}^2 f(w) dw - \int_0^{\hat{h}_k(x)} \left(1 - \frac{w}{\hat{h}_k(x)}\right)^2 f(w) dw \right| \\ & \leq (M + 2)(M + 1)F(\hat{h}_k(x)). \end{aligned}$$

With an application of (4.21), we know that

$$F(\hat{h}_k(x)) \leq F(\epsilon_k(x)) \leq \frac{Mx}{k\epsilon_k(x)} \leq \frac{M_F}{k},$$

since $x \leq \epsilon_k(x)$. As a result, if $\hat{h}_k(x) = x \leq \epsilon_k(x)$, then the left-hand side of (4.62) is bounded by $M_F(M+2)(M+1)/k$ from above.

Case (ii) $\hat{h}_k(x) = \epsilon_k(x) \leq x$. We apply Lemma 4.1 as well as the fundamental theorem of calculus to obtain that

$$\begin{aligned} & \left| \hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w) - \frac{w}{\hat{h}_k(x)} \right| = \left| \int_{x-w}^x \hat{v}'_{k-1}(y) - \frac{1}{\epsilon_k(x)} dy \right| \\ & \leq \int_{x-w}^x \frac{1}{\epsilon_k(y)} - \frac{1}{\epsilon_k(x)} dy + m \int_{x-w}^x \frac{1}{y^{3/4}\epsilon_k^{1/4}(y)} dy. \end{aligned} \quad (4.64)$$

The first integral in this last right-hand side has the closed-form representation

$$\int_{x-w}^x \frac{1}{\epsilon_k(y)} - \frac{1}{\epsilon_k(x)} dy = kF(\epsilon_k(x)) - kF(\epsilon_k(x-w)) - \frac{w}{\epsilon_k(x)},$$

which, by (4.23), is bounded from above by $kF(\epsilon_k(x))w^2/x^2$, that is,

$$\int_{x-w}^x \frac{1}{\epsilon_k(y)} - \frac{1}{\epsilon_k(x)} dy \leq \frac{kF(\epsilon_k(x))w^2}{x^2}. \quad (4.65)$$

The integrand of the second integral satisfies that

$$\frac{1}{y^{3/4}\epsilon_k^{1/4}(y)} = \frac{1}{4\epsilon_k(y)(kF(\epsilon_k(y)))^{3/4}} \cdot 4 \left(\frac{k\epsilon_k(y)F(\epsilon_k(y))}{y} \right)^{3/4}.$$

The first factor is the derivative of $(kF(\epsilon_k(y)))^{1/4}$, while the second factor is bounded by $4M_F^{3/4}$ due to (4.21). As a result, we have that,

$$\begin{aligned} & \int_{x-w}^x \frac{1}{y^{3/4}\epsilon_k^{1/4}(y)} dy \\ & \leq 4M_F^{3/4} \int_{x-w}^x d \left[(kF(\epsilon_k(y)))^{1/4} \right] = 4M_F^{3/4} \left\{ (kF(\epsilon_k(x)))^{1/4} - (kF(\epsilon_k(x-w)))^{1/4} \right\} \\ & = 4M_F^{3/4} \frac{kF(\epsilon_k(x)) - kF(\epsilon_k(x-w))}{\left\{ (kF(\epsilon_k(x)))^{1/4} + (kF(\epsilon_k(x-w)))^{1/4} \right\} \left\{ (kF(\epsilon_k(x)))^{1/2} + (kF(\epsilon_k(x-w)))^{1/2} \right\}} \\ & \leq 4M_F^{3/4} \frac{kF(\epsilon_k(x)) - kF(\epsilon_k(x-w))}{(kF(\epsilon_k(x)))^{3/4}}. \end{aligned}$$

If we apply (4.23) on the numerator of this last right-hand side, we obtain a further upper bound

$$\int_{x-w}^x \frac{1}{y^{3/4} \epsilon_k^{1/4}(y)} dy \leq 4M^{3/4} \frac{w/\epsilon_k(x) + kF(\epsilon_k(x))w^2/x^2}{(kF(\epsilon_k(x)))^{3/4}}.$$

By taking this last inequality and (4.65) back to (4.64), we obtain that

$$\begin{aligned} \left| \hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w) - \frac{w}{\hat{h}_k(x)} \right| &\leq \frac{kF(\epsilon_k(x))w^2}{x^2} + 4M^{3/4} \frac{w/\epsilon_k(x) + kF(\epsilon_k(x))w^2/x^2}{(kF(\epsilon_k(x)))^{3/4}} \\ &= 4M^{3/4} \frac{w/\epsilon_k(x)}{(kF(\epsilon_k(x)))^{3/4}} + \frac{kF(\epsilon_k(x))w^2}{x^2} \left\{ 1 + \frac{4M^{3/4}}{(kF(\epsilon_k(x)))^{3/4}} \right\}. \end{aligned}$$

Integrating this last equation against $dF(w)$ over $w \in [0, \epsilon_k(x)]$ leads us to

$$\begin{aligned} &\int_0^{\epsilon_k(x)} \left| \hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w) - \frac{w}{\hat{h}_k(x)} \right| f(w) dw \\ &\leq \frac{4M^{3/4}x}{k\epsilon_k(x)(kF(\epsilon_k(x)))^{3/4}} + \frac{\epsilon_k(x)F(\epsilon_k(x))}{x} \left\{ 1 + \frac{4M^{3/4}}{(kF(\epsilon_k(x)))^{3/4}} \right\} \\ &= \frac{4M^{3/4}F^{1/4}(\epsilon_k(x))}{k^{3/4}} \left\{ \frac{x}{k\epsilon_k(x)F(\epsilon_k(x))} + \frac{\epsilon_k(x)}{x} \right\} + \frac{\epsilon_k(x)F(\epsilon_k(x))}{x}. \end{aligned}$$

By (4.21), we know that the second summand is bounded by M_F/k . While for second factor of the first summand, we notice that $\epsilon_k(x)/x \leq 1$ when $\hat{h}_k(x) = \epsilon_k(x)$, and apply (4.21) again to obtain that

$$\frac{x}{k\epsilon_k(x)F(\epsilon_k(x))} \leq 1,$$

so overall we have the inequality

$$\int_0^{\epsilon_k(x)} \left| \hat{v}_{k-1}(x) - \hat{v}_{k-1}(x-w) - \frac{w}{\hat{h}_k(x)} \right| f(w) dw \leq \frac{8M^{3/4}F^{1/4}(\epsilon_k(x))}{k^{3/4}} + \frac{M_F}{k},$$

and the proof is complete. \square

Now, with this last proposition, we are ready to prove the key variance asymptotics result of this section.

Proof. (Proof of Lemma 4.2.) Recall the decomposition of $\omega_{n-t}(X_t)$ from (4.60)

$$\omega_{n-t}(X_t) = \sum_{j=t+1}^n \mathbb{E}[d_j^2 | \mathcal{F}_t] = \sum_{j=t+1}^n \mathbb{E}[B_j^2 | \mathcal{F}_t] - \sum_{j=t+1}^n \mathbb{E}[A_j^2 | \mathcal{F}_t],$$

and the two summands from (4.58) and (4.59) respectively

$$A_j^2 = \left(\int_0^{\hat{h}_{n-j+1}(X_{j-1})} \{1 + \hat{v}_{n-j}(X_{j-1} - w) - \hat{v}_{n-j}(X_{j-1})\} f(w) dw \right)^2$$

$$\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] = \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \{1 + \hat{v}_{n-j}(X_{j-1} - w) - \hat{v}_{n-j}(X_{j-1})\}^2 f(w) dw.$$

The remaining capacity process $\{X_t\}_{t=0}^n$ behaves nicely before the policy becomes greedy, so we split the time horizon $t = 0, 1, \dots, n$ into two parts: before and after the stopping time η , which is the first time that the policy $\hat{\pi}$ becomes greedy. Recall from (4.11) that the stopping time η is

$$\eta = \min \{t \in [n] : \epsilon_{n-t}(X_t) > X_t \text{ or } X_t > (n-t)\mu\}.$$

CONSERVATIVE HORIZON $1 \leq j \leq \eta$. The heuristic policy $\hat{\pi}$ uses the conservative threshold $\hat{h}_{n-j+1}(X_{j-1}) = \epsilon_{n-j+1}(X_{j-1})$. By Corollary 4.1 we know that the integrand of A_j is bounded by M , so

$$A_j^2 \leq M^2 F^2(\hat{h}_{n-j+1}(X_{j-1})),$$

and by (4.62) in Proposition 4.9, we know that

$$\left| \mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] - \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})}\right\}^2 f(w) dw \right|$$

$$\leq M \left(\frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}} + \frac{1}{n-j+1} \right).$$

Combine these last two estimations together, we obtain that

$$\left| \sum_{j=t+1}^{\eta} (\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=t+1}^{\eta} \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})}\right\}^2 f(w) dw \right|$$

$$\leq \sum_{j=t+1}^{\eta} \left\{ M \left(\frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}} + \frac{1}{n-j+1} \right) + M^2 F^2(\hat{h}_{n-j+1}(X_{j-1})) \right\}. \quad (4.66)$$

Since for $1 \leq j \leq \eta$, we have $\hat{h}_{n-j+1}(X_{j-1}) = \epsilon_{n-j+1}(X_{j-1}) \leq X_{j-1}$, we know that

$$\frac{F^2(\hat{h}_{n-j+1}(X_{j-1}))}{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))(n-j+1)^{1/4}}$$

$$= \left(\frac{(n-j+1)\epsilon_{n-j+1}(X_{j-1})F(\epsilon_{n-j+1}(X_{j-1}))}{X_{j-1}} \right)^{3/4} \left(\frac{F(\epsilon_{n-j+1}(X_{j-1}))}{\epsilon_{n-j+1}(X_{j-1})} \right)^{3/4} \left(X_{j-1}^{3/4} F^{1/4}(\epsilon_{n-j+1}(X_{j-1})) \right).$$

The first factor in this last right hand side is bounded from above by $M_F^{3/4}$ due to (4.21); the second factor is bounded from above by $m_u^{3/4}$ due to (4.12); the last factor is bounded by $c^{3/4}$ since $X_{j-1} \leq X_0 = c$ and the quantity $F(\epsilon_{n-j+1}(X_{j-1}))$ is always smaller than 1, so we obtain that

$$\frac{F^2(\hat{h}_{n-j+1}(X_{j-1}))}{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))(n-j+1)^{1/4}} \leq (M_F m_u c)^{3/4},$$

which in turn implies that

$$F^2(\hat{h}_{n-j+1}(X_{j-1})) \leq (M_F m_u c)^{3/4} \frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}}.$$

Substituting this inequality to the right-hand side of (4.66) leads us to

$$\begin{aligned} & \left| \sum_{j=t+1}^{\eta} (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=t+1}^{\eta} \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \\ & \leq \sum_{j=t+1}^{\eta} \left\{ (M + M^2(M_F m_u c)^{3/4}) \frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}} + \frac{M}{n-j+1} \right\} \\ & \leq (M + M^2(M_F m_u c)^{3/4}) \sum_{j=t+1}^{\tau} \frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}} + M(\log(n-t) + 1) \\ & \leq (M + M^2(M_F m_u c)^{3/4}) \sum_{j=t+1}^n \frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}} + M(\log(n-t) + 1). \end{aligned}$$

Take one more expectation on this last inequality with respect to the filtration \mathcal{F}_t and apply Proposition 4.3, we find that

$$\begin{aligned} & \mathbb{E} \left[\left| \sum_{j=t+1}^{\eta} (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=t+1}^{\eta} \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right] \\ & \leq 4(M + M^2(M_F m_u c)^{3/4}) \{(n-t)F(\epsilon_{n-t}(X_t))\}^{1/4} + M(\log(n-t) + 1) \end{aligned} \quad (4.67)$$

GREEDY HORIZON $\eta + 1 \leq j \leq n$. By the definition of η , we have either $X_\eta > (n - \eta)\mu$ or $X_\eta < \epsilon_{n-\eta}(X_\eta)$.

If $X_\eta > (n - \eta)\mu$, then we have that $n - \eta < X_\eta \mu^{-1} \leq X_0 \mu^{-1} = c\mu^{-1}$. As a result,

$$0 \leq \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \leq \sum_{j=\eta+1}^n F(\hat{h}_{n-j+1}(X_{j-1})) \leq n - \eta \leq c\mu^{-1}.$$

From the uniform boundedness of the martingale differences in (4.56), we also have the simple bound

$$\left\| \sum_{j=\eta+1}^n d_j^2 \right\|_{\infty} \leq 4(n-\eta)(M+1)^2 \leq 4c\mu^{-1}(M+1)^2,$$

so we find the almost sure bound

$$\left| \sum_{j=\eta+1}^n (\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \leq 4c\mu^{-1}(M+1)^2. \quad (4.68)$$

If $X_\eta < \epsilon_{n-\eta}(X_\eta)$, then we apply Corollary 4.1 on the integrand of $\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}]$ to obtain that

$$\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] \leq M^2 F(\hat{h}_{n-j+1}(X_{j-1})).$$

Then if we sum up this last estimation from $j = \eta + 1$ to n , we would get

$$0 \leq \sum_{j=\eta+1}^n \mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] - \sum_{j=\eta+1}^n A_j^2 \leq \sum_{j=\eta+1}^n \mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] \leq M^2 \sum_{j=\eta+1}^n F(\hat{h}_{n-j+1}(X_{j-1})).$$

Combining this last estimation with the simple bound

$$0 \leq \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \leq \sum_{j=\eta+1}^n F(\hat{h}_{n-j+1}(X_{j-1})),$$

we find that

$$\begin{aligned} & \left| \sum_{j=\eta+1}^n (\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \\ & \leq M^2 \sum_{j=\eta+1}^n F(\hat{h}_{n-j+1}(X_{j-1})). \end{aligned}$$

For any realization of η , if we take expectation of this last inequality conditioning on X_η , we obtain that

$$\begin{aligned} & \mathbb{E} \left[\left| \sum_{j=\eta+1}^n (\mathbb{E}[B_j^2 \mid \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \mid \eta, X_\eta \right] \\ & \leq M^2 \sum_{j=\eta+1}^n \mathbb{E} [F(\hat{h}_{n-j+1}(X_{j-1})) \mid X_\eta]. \end{aligned} \quad (4.69)$$

The inequality $X_\eta < \epsilon_{n-\eta}(X_\eta)$ tells us that the summation on this last right-hand side, which equals $\hat{v}_{n-\eta}(X_\eta)$, satisfies

$$\hat{v}_{n-\eta}(X_\eta) \leq (n-\eta)F(\epsilon_{n-\eta}(X_\eta)) = \frac{(n-\eta)\epsilon_{n-\eta}(X_\eta)F(\epsilon_{n-\eta}(X_\eta))}{X_\eta} \frac{X_\eta}{\epsilon_{n-\eta}(X_\eta)} \leq M_F, \quad (4.70)$$

due to (4.21). As a result, if we substitute this last upper bound to (4.69) and take one more expectation with respect to \mathcal{F}_t , we find that

$$\mathbb{E} \left[\left| \sum_{j=\eta+1}^n (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right] \leq M^2 M_F. \quad (4.71)$$

To sum up, in either case of (i) $X_\eta > (n-\eta)\mu$ or (ii) $X_\eta < \epsilon_{n-\eta}(X_\eta)$, from (4.68) and (4.71) we find that

$$\begin{aligned} & \mathbb{E} \left[\left| \sum_{j=\eta+1}^n (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right] \\ & \leq \max \{ M^2 M_F, 4c\mu^{-1}(M+1)^2 \}. \end{aligned} \quad (4.72)$$

COMBINING CONSERVATIVE AND GREEDY HORIZONS. From the conditional variance decomposition

$$\omega_{n-t}(X_t) = \sum_{j=t+1}^n \mathbb{E}[d_j^2 | \mathcal{F}_t] = \sum_{j=t+1}^n \mathbb{E}[B_j^2 | \mathcal{F}_t] - \sum_{j=t+1}^n \mathbb{E}[A_j^2 | \mathcal{F}_t],$$

if we combine (4.67) and (4.72) together and apply triangle inequality, we would get

$$\begin{aligned} & \mathbb{E} \left[\left| \sum_{j=t+1}^n (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=t+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right] \\ & \leq \mathbb{E} \left[\left| \sum_{j=t+1}^{\eta} (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=t+1}^{\eta} \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right] \\ & \quad + \mathbb{E} \left[\left| \sum_{j=\eta+1}^n (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=\eta+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right] \\ & \leq (M + M^2(M_F m_u c)^{3/4}) \{ (n-t)F(\epsilon_{n-t}(X_t)) \}^{1/4} + M(\log(n-t) + 1) + \max \{ M^2 M_F, 4c\mu^{-1}(M+1)^2 \} \end{aligned}$$

Note that we always have the relation

$$\begin{aligned} & \left| \omega_{n-t}(X_t) - \sum_{j=t+1}^n \mathbb{E} \left[\int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right] \right| \\ & \leq \mathbb{E} \left[\left| \sum_{j=t+1}^n (\mathbb{E}[B_j^2 | \mathcal{F}_{j-1}] - A_j^2) - \sum_{j=t+1}^n \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right| \middle| \mathcal{F}_t \right], \end{aligned}$$

so the proof will complete once we can show that

$$\begin{aligned} & \left| \gamma(r) \hat{v}_{n-t}(X_t) - \sum_{j=t+1}^n \mathbb{E} \left[\int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right] \right| \quad (4.73) \\ & \leq M \left\{ [(n-t)F(\epsilon_{n-t}(X_t))]^{1/4} + 1 \right\}. \end{aligned}$$

Recall that the value function $\hat{v}_{n-t}(X_t)$ satisfies that $\hat{v}_{n-t}(X_t) = \sum_{j=t+1}^n \mathbb{E} [F(\hat{h}_{n-j+1}(X_{j-1}))]$,

so the left-hand side of this last inequality can be represented as

$$\left| \sum_{j=t+1}^n \mathbb{E} \left[\gamma(r) F(\hat{h}_{n-j+1}(X_{j-1})) - \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right] \right|.$$

To analyze each summand in this last representation, we introduce the shorthand function $\eta_k(x) : [0, c] \rightarrow \mathbb{R}$ defined as

$$\eta_k(x) \equiv \gamma(r) F(\epsilon_k(x)) - \int_0^{\epsilon_k(x)} \left\{ 1 - \frac{w}{\epsilon_k(x)} \right\}^2 f(w) dw,$$

and consider the quantity $k^{3/4} \eta_k(x) / F^{1/4}(\epsilon_k(x))$. In fact, we have the identity

$$\begin{aligned} & \frac{k^{3/4} \eta_k(x)}{F^{1/4}(\epsilon_k(x))} \\ &= \frac{k^{3/4} \int_0^{\epsilon_k(x)} \left[\frac{2r^2}{(1+r)(1+2r)} - \left(1 - \frac{w}{\epsilon_k(x)} \right)^2 \right] f(w) dw}{F^{1/4}(\epsilon_k(x))} \\ &= \frac{2k^{3/4} \int_0^{\epsilon_k(x)} \left(\frac{w}{\epsilon_k(x)} - \frac{1}{1+r} \right) f(w) dw}{F^{1/4}(\epsilon_k(x))} + \frac{k^{3/4} \int_0^{\epsilon_k(x)} \left(\frac{1}{1+2r} - \frac{w^2}{\epsilon_k^2(x)} \right) f(w) dw}{F^{1/4}(\epsilon_k(x))} \\ &= 2k^{3/4} F^{3/4}(\epsilon_k(x)) \left(\frac{r}{1+r} - \frac{\int_0^{\epsilon_k(x)} F(w) dw}{\epsilon_k(x) F(\epsilon_k(x))} \right) + 2k^{3/4} F^{3/4}(\epsilon_k(x)) \left(\frac{\int_0^{\epsilon_k(x)} w F(w) dw}{\epsilon_k^2(x) F(\epsilon_k(x))} - \frac{r}{1+2r} \right) \end{aligned}$$

Now, if we take absolute value on both sides and apply triangle inequality, we obtain that

$$\frac{k^{3/4} |\eta_k(x)|}{F^{1/4}(\epsilon_k(x))} \leq 2k^{3/4} F^{3/4}(\epsilon_k(x)) \left\{ \left| \frac{r}{1+r} - \frac{\int_0^{\epsilon_k(x)} F(w) dw}{\epsilon_k(x) F(\epsilon_k(x))} \right| + \left| \frac{\int_0^{\epsilon_k(x)} w F(w) dw}{\epsilon_k^2(x) F(\epsilon_k(x))} - \frac{r}{1+2r} \right| \right\}.$$

By (4.19), we have the further bound

$$\frac{k^{3/4} |\eta_k(x)|}{F^{1/4}(\epsilon_k(x))} \leq 2M (k\epsilon_k(x) F(\epsilon_k(x)))^{3/4},$$

where this last right-hand side is bounded by $2M(M_F x)^{3/4}$ due to (4.21). As a result,

$$|\eta_k(x)| \leq 2M(M_F x)^{3/4} \frac{F^{1/4}(\epsilon_k(x))}{k^{3/4}} \leq 2M(M_F c)^{3/4} \frac{F^{1/4}(\epsilon_k(x))}{k^{3/4}}.$$

This last estimation on $\eta_k(x)$ tells us that

$$\begin{aligned} & \left| \sum_{j=t+1}^{\eta} \mathbb{E} \left[\gamma(r) F(\hat{h}_{n-j+1}(X_{j-1})) - \int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \mid \mathcal{F}_t \right] \right| \\ & \leq \sum_{j=t+1}^{\eta} \mathbb{E} [|\eta_{n-j+1}(X_{j-1})| \mid \mathcal{F}_t] \leq 2M(M_F c)^{3/4} \sum_{j=t+1}^{\eta} \mathbb{E} \left[\frac{F^{1/4}(\epsilon_{n-j+1}(X_{j-1}))}{(n-j+1)^{3/4}} \mid \mathcal{F}_t \right]. \end{aligned}$$

An application of Proposition 4.3 yields that this last right-hand side has a further upper bound $8M(M_F c)^{3/4} \{(n-t)F(\epsilon_{n-t}(X_t))\}^{1/4}$. Lastly, (4.73) will become evident once one notices that the uncertainties after the stopping time τ is negligible. That is, $\hat{v}_{n-\tau}(X_\tau) \leq M_F$ a.s., and

$$\begin{aligned} \sum_{j=\eta+1}^n \mathbb{E} \left[\int_0^{\hat{h}_{n-j+1}(X_{j-1})} \left\{ 1 - \frac{w}{\hat{h}_{n-j+1}(X_{j-1})} \right\}^2 f(w) dw \right] & \leq \sum_{j=\eta+1}^n \mathbb{E} [F(\hat{h}_{n-j+1}(X_{j-1}))] \\ & \leq \max\{M_F, K_0\} \quad \text{a.s.,} \end{aligned}$$

due to (4.70) and the definition of η . □

4.4 Martingale Central Limit Theorem

Our proof of Theorem 4.1 relies on a basic version of the martingale central limit theorem which is contained in, e.g., McLeish (1974, Theorem 2.3).

Proposition 4.10 (Martingale central limit theorem). *Let $\{Z_j\}_{j=1}^n$ be a martingale difference sequence with respect to the increasing σ -field $\{\mathcal{F}_j\}_{j=1}^n$. If*

$$\max_{1 \leq j \leq n} \|Z_j\|_\infty \rightarrow 0 \quad \text{as } n \rightarrow \infty, \quad (4.74)$$

and

$$\sum_{j=1}^n \mathbb{E}[Z_j^2 \mid \mathcal{F}_{j-1}] \xrightarrow{p} 1 \quad \text{as } n \rightarrow \infty, \quad (4.75)$$

then we have the convergence in distribution

$$\sum_{j=1}^n Z_j \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty.$$

In order to verify (4.75) in our context, we need the following proposition.

Proposition 4.11 (A second conditional variance bound). *There exists a constant M that depends only on c and F such that for all $n \geq 1$,*

$$\text{Var} \left\{ \sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \right\} \leq M(nF(\epsilon_n(c))) \left([nF(\epsilon_n(c))]^{1/2} + (\log n)^2 + 1 \right).$$

Proof. To simplify notations, we let $V = \sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}]$ and define another martingale as

$$V_t \equiv \mathbb{E}[V \mid \mathcal{F}_t] \quad \text{for all } t \in [n].$$

Then we define the associated martingale differences as $\Delta_t \equiv V_t - V_{t-1}$ for all $t \in [n]$. For each $1 \leq t \leq n$, since $\sum_{j=1}^t \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}]$ is \mathcal{F}_{t-1} -measurable, by the definition of Δ_t we have

$$\begin{aligned} \Delta_t &= \mathbb{E} \left[\sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \mid \mathcal{F}_t \right] - \mathbb{E} \left[\sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \mid \mathcal{F}_{t-1} \right] \\ &= \mathbb{E} \left[\sum_{j=t+1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \mid \mathcal{F}_t \right] - \mathbb{E} \left[\sum_{j=t+1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \mid \mathcal{F}_{t-1} \right] \\ &= \sum_{j=t+1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_t] - \mathbb{E} \left[\sum_{j=t+1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_t] \mid \mathcal{F}_{t-1} \right] \\ &= \omega_{n-t}(X_t) - \mathbb{E}[\omega_{n-t}(X_t) \mid \mathcal{F}_{t-1}]. \end{aligned}$$

Recall the event of selecting the t -th item

$$G_t \equiv \left\{ \omega : W_t(\omega) \in [0, \hat{h}_{n-t+1}(X_{t-1}(\omega))] \right\},$$

and condition on the occurrence of G_t , we have that

$$\omega_{n-t}(X_t) = \omega_{n-t}(X_{t-1}) + (\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1})) \mathbb{1}\{G_t\}.$$

Substitute this back to the representation of Δ_t , we find that

$$\Delta_t = (\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1})) \mathbb{1}\{G_t\} - \mathbb{E}[(\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1})) \mathbb{1}\{G_t\} \mid \mathcal{F}_{t-1}].$$

Hence by squaring this last identity and taking expectation with respect to the filtration \mathcal{F}_{t-1} , we obtain the conditional second moment bound

$$\begin{aligned} \mathbb{E}[\Delta_t^2 \mid \mathcal{F}_{t-1}] &= \mathbb{E}[(\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1}))^2 \mathbb{1}\{G_t\} \mid \mathcal{F}_{t-1}] \\ &\quad - \mathbb{E}[(\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1})) \mathbb{1}\{G_t\} \mid \mathcal{F}_{t-1}]^2 \\ &\leq \mathbb{E}[(\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1}))^2 \mathbb{1}\{G_t\} \mid \mathcal{F}_{t-1}]. \end{aligned}$$

By Lemma 4.2 and the triangle inequality, we know that

$$\begin{aligned} &|\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1})| \mathbb{1}\{G_t\} \\ &\leq |\omega_{n-t}(X_{t-1} - W_t) - \gamma(r)\hat{v}_{n-t}(X_{t-1} - W_t)| \mathbb{1}\{G_t\} + |\omega_{n-t}(X_{t-1}) - \gamma(r)\hat{v}_{n-t}(X_{t-1})| \\ &\quad + \gamma(r)|\hat{v}_{n-t}(X_{t-1} - W_t) - \hat{v}_{n-t}(X_{t-1})| \mathbb{1}\{G_t\} \\ &\leq 2M \left\{ [(n-t)F(\epsilon_{n-t}(X_{t-1}))]^{1/4} + \log(n-t) + 1 \right\} + \gamma(r)|\hat{v}_{n-t}(X_{t-1} - W_t) - \hat{v}_{n-t}(X_{t-1})| \mathbb{1}\{G_t\}. \end{aligned}$$

Corollary 4.1 tells us that the second summand in this last right-hand side is bounded by $M\gamma(r)$, so

$$|\omega_{n-t}(X_{t-1} - W_t) - \omega_{n-t}(X_{t-1})| \mathbb{1}\{G_t\} \leq 2M \left\{ [(n-t)F(\epsilon_{n-t}(X_{t-1}))]^{1/4} + \log(n-t) + 1 \right\} + M\gamma(r),$$

and hence

$$\mathbb{E}[\Delta_t^2 \mid \mathcal{F}_{t-1}] \leq M \left\{ (nF(\epsilon_n(c)))^{1/2} + (\log n)^2 + 1 \right\} \mathbb{E}[\mathbb{1}\{G_t\} \mid \mathcal{F}_{t-1}]$$

As a result, we find that

$$\begin{aligned} \text{Var} \left\{ \sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \right\} &= \sum_{t=1}^n \mathbb{E}[\Delta_t^2] = \sum_{t=1}^n \mathbb{E}[\mathbb{E}[\Delta_t^2 \mid \mathcal{F}_{t-1}]] \\ &\leq M \left\{ (nF(\epsilon_n(c)))^{1/2} + (\log n)^2 + 1 \right\} \sum_{t=1}^n \mathbb{E}[\mathbb{E}[\mathbb{1}\{G_t\} \mid \mathcal{F}_{t-1}]] \\ &= M \left\{ (nF(\epsilon_n(c)))^{1/2} + (\log n)^2 + 1 \right\} \mathbb{E}[\hat{N}_n(c)] \leq M (nF(\epsilon_n(c))) \left\{ (nF(\epsilon_n(c)))^{1/2} + (\log n)^2 + 1 \right\}, \end{aligned}$$

since $\mathbb{E}[\hat{N}_n(c)] \leq nF(\epsilon_n(c))$. The proof is complete. \square

Now, we have all the ingredients to prove Theorem 4.1.

Proof. (Proof of Theorem 4.1.) If we define the normalized martingale differences as

$$Z_j \equiv \frac{d_j}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \quad \text{for all } 1 \leq j \leq n,$$

then the uniform boundedness (4.56) of the martingale differences together with the growth condition (4.16) imply the negligibility condition (4.74).

With the growth condition (4.16), the conditional variance bound (4.61) tells us that

$$\sum_{j=1}^n \mathbb{E}[d_j^2] = \omega_n(X_0) \sim \gamma(r)nF(\epsilon_n(c)) \quad \text{as } n \rightarrow \infty.$$

On the other hand, we know from Proposition 4.11 that

$$\text{Var} \left\{ \sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \right\} \sim o([nF(\epsilon_n(c))]^2) \quad \text{as } n \rightarrow \infty.$$

Hence if we combine these last two asymptotics and apply Chebyshev's inequality, we obtain that for any $\delta > 0$ fixed,

$$\begin{aligned} \mathbb{P} \left\{ \left| \sum_{j=1}^n \mathbb{E}[Z_j^2 \mid \mathcal{F}_{j-1}] - 1 \right| > \delta \right\} &= \mathbb{P} \left\{ \left| \sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] - \gamma(r)nF(\epsilon_n(c)) \right| > \gamma(r)nF(\epsilon_n(c))\delta \right\} \\ &\leq \frac{\text{Var} \left\{ \sum_{j=1}^n \mathbb{E}[d_j^2 \mid \mathcal{F}_{j-1}] \right\}}{[nF(\epsilon_n(c))]^2 \gamma(r)^2 \delta^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Therefore, the convergence in probability condition (4.75) is satisfied and so by Proposition 4.10 we have the convergence in distribution

$$\sum_{j=1}^n Z_j = \sum_{j=1}^n \frac{d_j}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} = \frac{\hat{N}_n(c) - \mathbb{E}[\hat{N}_n(c)]}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty.$$

Finally, the identity $\mathbb{E}[\hat{N}_n(c)] = nF(\epsilon_n(c)) + O(\log n)$ from (4.15) together with the growth condition (4.16) lead us to the weak law

$$\frac{\hat{N}_n(c) - nF(\epsilon_n(c))}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty.$$

For the weak convergence of $N_n^*(c)$, we apply Theorem 2.1 to obtain that

$$\frac{\mathbb{E} \left[\left| N_n^*(c) - \hat{N}_n(c) \right| \right]}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} = \frac{O(\log n)}{\sqrt{\gamma(r)nF(\epsilon_n(c))}}.$$

The growth condition (4.16) then implies that the last right-hand side vanishes as $n \rightarrow \infty$. As a result,

$$\frac{N_n^*(c) - \hat{N}_n(c)}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \rightarrow 0 \quad \text{in } L^1,$$

and hence we have that

$$\frac{N_n^*(c) - nF(\epsilon_n(c))}{\sqrt{\gamma(r)nF(\epsilon_n(c))}} \Rightarrow N(0, 1) \quad \text{as } n \rightarrow \infty,$$

just as we wanted. □

4.5 Concluding Remarks

In this chapter, we study the dynamic and stochastic knapsack problem with unitary rewards and revisit the adaptive heuristic policy $\hat{\pi}$ of Chapter 2. We show that the number of selections made by such heuristic shares the same variance asymptotics and limiting distribution with that made the optimal dynamic programming policy. However, as we have noted, different asymptotically optimal policies may have different higher moments as well as different limiting behaviors. Such differences may be overlooked if asymptotic optimality is the only criterion to judge different policies.

Looking forward, we believe that the analysis of higher moments and the limiting distribution of the total reward collected by different policies is important when studying sequential decision problems. The analysis in this chapter serves as a first step towards this direction, and other related problems could potential by studied in a similar spirit. In addition to the mean-performance and worst-case-performance analyses, the complementary analysis of the higher moments and limiting behaviors of the total rewards collected by different policies would provide valuable and comprehensive managerial insights to decision makers.

Chapter 5

Data-driven Monitoring the Implementation of Policies across Many Markov Decision Problems

Markov decision problem (or MDP), with both finite and infinite horizon, is a commonly used framework for modeling sequential decision problems under uncertainty, and have been extensively studied in the literature. Numerous applications of MDP theory can be found in finance, operations management, operations research, robotics and telecommunication, just to name a few (see, e.g. White, 1993; Puterman, 1994; Bertsekas, 2005). In this chapter, we focus our attention to many discrete time finite-horizon MDPs faced in parallel by self-interested decision makers, where the exogenous uncertainties faced by each decision maker have arbitrary and unknown correlation structure. Our goal is to simultaneously monitor whether each of the decision makers is implementing her designated policy, with only the knowledge of the marginal distributions of the exogenous uncertainty but not the correlation structure. The key feature of our monitoring is that the number of decision makers, p , is allowed to be much larger than the number of periods, n . We also provide data-driven procedure to construct simultaneous confidence intervals for the total rewards collected by each decision maker, with theoretical asymptotic validity.

In a typical finite horizon MDP with n discrete time periods, the mean-optimal policy π^* is the policy that delivers the maximal expected total reward. That is,

$$\pi^* = \arg \max_{\pi \in \Pi(n)} \mathbb{E}[R_n(\pi)],$$

where $\Pi(n)$ is the class of all non-anticipating sequential policies. It is well-known that under mild conditions, there exists a Markovian deterministic policy that is optimal. Such optimal policy π^* is often defined through Bellman equation with the associated value function given by

$$v_n^* = \mathbb{E}[R_n(\pi^*)] = \sup_{\pi \in \Pi(n)} \mathbb{E}[R_n(\pi)],$$

However, the total reward $R_n(\pi^*)$ collected by the optimal policy π^* is *ex ante* random, and its expected value v_n^* alone does not tell us about where $R_n(\pi^*)$ would actually land. The same

This chapter is written under the supervision of Prof. Alexandre Belloni. The results presented here are also in a joint research paper Belloni and Xie (2020).

randomness exists for the total reward collected by other near-optimal policies as well. When there are p correlated MDPs faced in parallel by self-interested decision makers labeled by $j \in [p] \equiv \{1, 2, \dots, p\}$, the vector of total rewards $\mathbf{R}_n \in \mathbb{R}^p$ each of them collects is also random, and it undergoes fluctuations and correlations that are generated both endogenously (sequential decisions in response to uncertainty realizations) and exogenously (correlation of the exogenous uncertainty). Without knowing the correlation structure of the exogenous uncertainty, but only knowing their marginal distributions instead, we are interested in testing the *null hypothesis*

$$H_0 \equiv \{ \text{every decision maker } j \in [p] \text{ is implementing her designated policy } \pi_j \}. \quad (5.1)$$

That is, we are interested in simultaneously monitoring many individual decision makers with the observation of the rewards they collect. We are also interested in constructing simultaneous confidence intervals for the total rewards $\mathbf{R}_n \in \mathbb{R}^p$ collected by each individual decision maker, with asymptotic validity.

For the case of one MDP, there are rich literature on limit theorems for the total reward collected by the optimal policy. Although the majority work along this line of research has focused on the infinite-horizon formulation in which the optimal policy is stationary, a notable exception is Arlotto and Steele (2016) who focuses on a finite-horizon formulation. They generalize the work of Dobrushin (1956a,b) on additive functionals of non-stationary Markov chains, and show that, in the finite-horizon MDP setting, under proper condition on the Dobrushin's contraction coefficient, the total reward collected by the optimal policy satisfies a central limit theorem (CLT) as the number of time periods, n , goes to infinity.

When CLTs are available, the above-mentioned hypothesis testing and the construction of confidence intervals are just simple exercises. However, when the “dimension” p becomes large, the CLTs start to fail, even for partial sums of independent and identically distributed (i.i.d.) random variables. In fact, for i.i.d. zero-mean random variables $Y_1, \dots, Y_n \in \mathbb{R}^p$ with $\mathbb{E}[Y_1 Y_1'] = I_p$ (p -by- p identity matrix), and their Gaussian counterparts $\tilde{Y}_i \sim N(0, I_p)$ for $i \in [n]$, Nagaev (1976) proves that

$$\sup_{A \in \mathcal{A}} \left| \mathbb{P} \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i \in A \right\} - \mathbb{P} \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{Y}_i \in A \right\} \right| \geq c \frac{\mathbb{E}[\|Y_1\|^3]}{\sqrt{n}}, \quad (5.2)$$

where \mathcal{A} is the class of all Borel measurable convex sets and $c > 0$ is some universal constant. This

last result, together with an simple application of Jensen's inequality

$$\mathbb{E} \left[\|Y_1\|^3 \right] \geq \left(\mathbb{E} \left[\|Y_1\|^2 \right] \right)^{3/2} = p^{3/2},$$

tells us that when $p \sim n^{1/3}$, the right-hand side of (5.2) is strictly bounded away from zero in the limiting regime $n \rightarrow \infty$, and hence a CLT simply cannot hold. Nevertheless, the groundbreaking work of Chernozhukov et al. (2013) shows that, instead of (5.2), one has that

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left\{ \max_{j \in [p]} \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_{ij} \leq t \right\} - \mathbb{P} \left\{ \max_{j \in [p]} \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{Y}_{ij} \leq t \right\} \right| \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

where p is allowed to be exponentially larger than n . This powerful result, which is very suitable for simultaneous hypothesis testing in high-dimensional regime, is then generalized and utilized by several papers to high-dimensional hypothesis testing problem under different weak dependence structures. Specifically, Zhang and Cheng (2014), Zhang and Wu (2017), and Zhang and Cheng (2018) focus on weakly dependent stationary time series that are characterized in terms of functional dependence, a measure of dependency introduced in Wu (2005); stationary β -mixing (absolute regular) sequence of random vectors is treated in the extension of Chernozhukov et al. (2019); while Belloni and Oliveira (2018) focuses on sequence of martingale differences.

Contribution and Organization of the Chapter

The main contribution of this chapter is to bridge the line of research in high-dimensional statistics and the theory of finite horizon MDPs, where dependency and non-stationarity are endogenously generated by decision makers. The rest of the chapter is organized as follows. We review related literature in Section 5.1. In Section 5.2, we introduce a general MDP framework that we will be working with. Then we take up the problems of hypothesis testing and the construction of simultaneous confidence intervals. Such problems are then treated in Section 5.3 and Section 5.4. Roughly speaking, Section 5.3 deals with policies that make states of MDPs “regenerative”; while Section 5.4 deals with MDPs with “transient” states. Then in Section 5.5, we consider a dynamic inventory control problem to which we apply the results of Section 5.3. Later in Section 5.6, we revisit the problem of sequential monotone subsequence selection from a random sample, that fits in the results of Section 5.4. Finally, we conclude in Section 5.7.

5.1 Literature Review

Markov decision problems, which also bear other names such as stochastic dynamic program, sequential stochastic optimization, and discrete time stochastic control, have been extensively studied in the literature. Many structural results have been obtained for the infinite-horizon formulation of such problems, in which the objective is usually the expected total reward with discount or the expected long-run average reward. In such formulations, the optimal policy is often stationary (time-independent). Asymptotic behaviors of the (random) total reward collected under stationary policies have been studied in great details in the literature. Substantial results have been attained through the works of Mandl and Laušmanová (1991); Mendoza-Pérez (2008); Mendoza-Pérez and Hernández-Lerma (2010). On the other hand, however, the finite-horizon case is more delicate and much less touched. A recent work that deserves attention is Arlotto and Steele (2016), who study a finite-horizon MDP, and establish the limiting distribution of the total reward collected under the optimal policy that is non-stationary (time-dependent).

The total reward collected under any sequential policy of a Markov decision problem can be viewed as the sum of the reward function on the states of the MDP. Since the states of the MDP form a Markov chain whose transition law is determined by the policy, the analysis of the total collected reward can also be analyzed through *Markov chain central limit theorems*. Maxwell and Woodroffe (2000) established asymptotic normality for partial sum of additive functionals on stationary ergodic Markov chain under moment conditions. Drift conditions are also proposed to guarantee central limit theorems on Markov chains, see, e.g. Jarner and Roberts (2002, Theorem 4.2) and Meyn and Tweedie (2009, Theorem 17.0.1). Again, a vast majority of this literature also focus on the stationary case. An excellent survey and extensive references on Markov chain CLT can be found in Jones (2004).

In addition to Markov chain central limit theorems, *mixing conditions* have also been widely used for establishing limiting results. The applications of such conditions include the context of Markov chains, as well as dependent sequences and time series. Most attention along this stream of literature has been put on the stationary case. Peligrad (1985, 1990) proves (functional) central limit theorems for stationary φ -mixing sequences under different moment conditions. Utev (1991) studies φ -mixing sequence and establishes a central limit theorem under a Lindeberg type condition. An exception that deserves special attention is Peligrad (2012), who proves a central limit theorem

for the sum of additive functionals on a non-homogeneous Markov chain whose ρ -mixing coefficient satisfies $\rho_1 < 1$. Another work for the non-stationary case is Fryzlewicz and Rao (2011), in which the mixing rates of time-varying ARCH time series models were obtained. An excellent review with more sustained references about different mixing conditions and their implications can be found in Bradley (2005).

All the above mentioned literature have been focused on proving limiting results in the classic setup where the dimension is fixed. However, as we discussed earlier, certain types of statistical inferences are still possible even when we are in the high-dimensional regime in which such above-mentioned results fail. This high-dimensional regime has drawn significant attention within the statistics literature. A series of remarkable results for the independent case have been achieved along the works of Chernozhukov et al. (2013, 2015, 2017). With these results, the applicability of Gaussian approximations have been pushed much further to the regime where the dimension of the estimand could be much larger than the sample size. See also the recent survey Belloni et al. (2018). Another notable recent work that deserves special attention is Belloni and Oliveira (2018), who develop high-dimensional central limit theorems and simultaneous inferences under martingale settings.

Different from the situation in low-dimension regime, the estimation of the covariance matrices under dependency is cumbersome when the dimension is large. As a result, the classic statistical inference approaches are not applicable any more. As a non-parametric procedure, *bootstrap* and its different variants become useful for high-dimensional inferences. In Chernozhukov et al. (2019), together with the block multiplier bootstrap method, β -mixing condition is used to construct simultaneous confidence intervals in high dimension. This particular bootstrap method can be thought of as a special version of the tapered bootstrap that has been studied in Paparoditis and Politis (2001, 2002) and Andrews (2004). We also refer interested readers to Lahiri (2003) as a general reference about resampling methods for dependent data. Different than β -mixing, Wu (2005) introduces another notion of dependence called functional dependence, and it is subsequently used in Zhang and Cheng (2014); Zhang and Wu (2017) and Zhang and Cheng (2018) to establish the validity of bootstrapping for stationary time series data in high dimension. We note that neither of the two weak-dependency conditions implies the other.

5.2 A General MDP Framework

We consider in total p decision makers, labeled by $j \in [p]$, each of whom faces a discrete-time finite horizon MDP with n periods, labeled by $t \in [n]$. Denote \mathcal{X}_j as the *state space* for the j -th decision maker (referred to as *she*). During each period, she takes an action that will have impact on the future rewards she collects. To incorporate different applications, we assume that there are uncertainties realized both before and after each action is taken. Specifically, at the beginning of each time period $t \in [n]$, suppose her MDP is in state $x_{(t-1)j}$, there first realizes an exogenous uncertainty $y_{tj} \in \mathcal{Y}_j$, then she chooses an action $a_{tj} \in \mathcal{A}(x_{(t-1)j}, y_{tj})$. The policy π_j that she implements is represented through how such action is chosen. Here one can think of the policy π_j as the optimal policy or a near-optimal policy. After such action is taken, another uncertainty $w_{tj} \in \mathcal{W}_j$ is realized, then she obtains a reward $r_j(x_{(t-1)j}, y_{tj}, a_{tj}, w_{tj})$, and the state of her MDP involves as

$$x_{tj} = f_j(x_{(t-1)j}, y_{tj}, a_{tj}, w_{tj}),$$

where r_j and f_j are two deterministic functions. The total reward she collects under policy π_j is

$$R_{nj} = R_{nj}(\pi_j) = \sum_{t=1}^n r_j(X_{(t-1)j}, Y_{tj}, A_{tj}, W_{tj}). \quad (5.3)$$

Denote $v_{nj}(x)$ as the expected total reward she collects by implementing policy π_j when the initial state of her MDP is $X_{0j} = x$. That is,

$$v_{nj}(x) = \mathbb{E}[R_{nj}(\pi_j)].$$

Similarly, we define her expected reward-to-go $v_{(n-k)j}(x)$ as the expected total reward that she collects when there are $n - k$ periods remaining and the state of her MDP is $X_{kj} = x$. That is,

$$v_{(n-k)j}(x) = \mathbb{E} \left[\sum_{t=k+1}^n r_j(X_{(t-1)j}, Y_{tj}, A_{tj}^{\pi_j}, W_{tj}^{\pi_j}) \mid X_{kj} = x \right], \quad (5.4)$$

where the notations A^{π_j} and W^{π_j} indicate that they are under policy π_j .

Sklar's Theorem and Copulas

We assume that all exogenous uncertainties faced by each decision maker have joint distribution F , and such uncertainties are independent across time periods. For any multi-dimensional distribution

F with marginal distributions $\{F_j\}_{j \in [p]}$, Sklar's theorem (Sklar, 1959) tells us that the joint distribution can be equivalently represented as a copula function of all the marginals. Specifically, one has that

$$F(w_1, \dots, w_p) = C(F_1(w_1), \dots, F_p(w_p)) \quad \text{for all } (w_1, \dots, w_p) \in \mathbb{R}^p.$$

The copula $C(\cdot) : [0, 1]^p \rightarrow [0, 1]$ together with all the marginal distributions $\{F_j\}_{j \in [p]}$ totally determine the entire distribution F . However, an arbitrary unknown copula C could put complex unknown correlation structure on the total rewards collected by each decision maker, even if all the marginal distributions are known. In other words, assuming knowing marginal distributions of the uncertainty puts the problem into a non-trivial yet tractable context.

5.2.1 Hypothesis Testing and the Construction of Simultaneous Confidence Intervals

We are interested in testing the null hypothesis

$$H_0 = \{ \text{every decision maker } j \in [p] \text{ is implementing her designated policy } \pi_j \},$$

where the number of decision makers, p , is allowed to be comparable to or even much larger than the number of time periods, n . We are also interested in constructing simultaneous confidence intervals for the total rewards collected by each decision maker under the null hypothesis. That is, for the vector

$$\mathbf{R}_n \equiv (R_{n1}(\pi_1), R_{n2}(\pi_2), \dots, R_{np}(\pi_p))' \in \mathbb{R}^p,$$

and prescribed confidence level α , we intend to construct simultaneous confidence intervals $\prod_{j=1}^p I_j(n, \alpha)$, with $I_j(n, \alpha) \subset \mathbb{R}$ for each $j \in [p]$, such that under the null hypothesis H_0 ,

$$\mathbb{P}\{R_{nj} \in I_j(n, \alpha) \text{ for all } j \in [p]\} \rightarrow 1 - \alpha, \quad \text{as } n \rightarrow \infty.$$

5.2.2 Martingale Representation

For the total reward collected by a policy of a finite horizon MDPs, one can define a Doob's martingale with respect to the filtration $\{\mathcal{F}_t\}_{t \in [n]}$, which in turn induces a sequence of martingale differences. Specifically, let \mathcal{F}_0 be the trivial σ -field and let $\mathcal{F}_t = \sigma\{Y_1, W_1, \dots, Y_t, W_t\}$ for $t \in [n]$, then for all $t \in [n]$,

$$\mathbf{d}_t \equiv \mathbb{E}[\mathbf{R}_n \mid \mathcal{F}_t] - \mathbb{E}[\mathbf{R}_n \mid \mathcal{F}_{t-1}] \in \mathbb{R}^p$$

compose a sequence of martingale differences in \mathbb{R}^p that satisfy that $\mathbb{E}[\mathbf{d}_t \mid \mathcal{F}_{t-1}] = 0$ for all $t \in [n]$ and

$$\sum_{t=1}^n \mathbf{d}_t = \mathbf{R}_n - \mathbb{E}[\mathbf{R}_n]. \quad (5.5)$$

Since each martingale difference \mathbf{d}_t is \mathcal{F}_t -measurable, from the total reward representation (5.3) and the value function representation (5.4), we have another useful representation that for the j -th coordinate of \mathbf{d}_t ,

$$d_{tj} = v_{(n-t)j}(X_{tj}) + r_j(X_{tj}, Y_{tj}, A_{tj}, W_{tj}) - v_{(n-t+1)j}(X_{(t-1)j}) \quad \text{for all } j \in [p]. \quad (5.6)$$

We close this section with the observation that with this last representation, we can construct test statistics for our statistical inference problems. Denote \hat{d}_{tj} as the realization of the martingale difference d_{tj} . That is, if we let the tuple $(x_{(t-1)j}, y_{tj}, a_{tj}, w_{tj})$ be the *realization* of $(X_{(t-1)j}, Y_{tj}, A_{tj}, W_{tj})$, then the state x_{tj} is determined as

$$x_{tj} = f_j(x_{(t-1)j}, y_{tj}, a_{tj}, w_{tj}),$$

and so we have the realized (observed) martingale difference

$$\hat{d}_{tj} = v_{(n-t)j}(x_{tj}) + r_j(x_{(t-1)j}, y_{tj}, a_{tj}, w_{tj}) - v_{(n-t+1)j}(x_{(t-1)j}). \quad (5.7)$$

The sequence of the realized martingale differences will serve as the building block of our bootstrap procedure later.

5.3 MDPs with Regenerative Property

In certain types of MDPs, the system states under certain policies enjoy regenerations in the sense that with high probability, the system state becomes independent of past history that is far away. In other words, the dependency between the states at two distant time periods decays as the distance gets large. When the speed of such decay is well balanced with the dimension p , we are able to conduct a block version of multiplier bootstrap method to obtain critical value for simultaneous hypothesis testing and to construct simultaneous confidence intervals. The decaying condition we need is made precise in the following definition.

Definition 5.1 (Tail of regeneration time.). Let $\{\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_n\} \subset \prod_{j \in [p]} \mathcal{X}_j$ be the sequence of (joint) states of each decision maker, that is,

$$\mathbf{X}_t = (X_{t1}, X_{t2}, \dots, X_{tp}) \in \prod_{j \in [p]} \mathcal{X}_j \quad \text{for all } t = 0, 1, \dots, n.$$

We say that the MDPs satisfy regenerative property with rate $\varphi(p, k)$ if for each t , there exists a random regeneration time τ such that $\mathbf{X}_{\tau+t}$ is independent of \mathbf{X}_t , and the tail probability of τ satisfies $\mathbb{P}\{\tau > k\} \leq \varphi(p, k)$.

The interplay between p and k in the decay rate function $\varphi(p, k)$ is the driving force of how large p is allowed for our statistical inferences. We close this part with the following remark before diving into the multiplier bootstrap procedure for testing our null hypothesis.

Remark 5.1. (φ -mixing vs. β -mixing.) If the function $\varphi(p, k) \rightarrow 0$ as $k \rightarrow \infty$, then the sequence $\{\mathbf{X}_t\}_{t \in [n]}$ is φ -mixing (see, e.g., Fan and Yao, 2003, Section 2.6). In fact, this mixing condition can be further relaxed to β -mixing. However, β -mixing condition is usually difficult to verify in applications, while on the other hand, the regeneration condition is common and much easier to verify in certain concrete MDPs.

5.3.1 Test Statistic and the Block Multiplier Bootstrap

Often times, for MDP policies that induce regenerative states, the variance of the total reward they collect grows linearly in terms of time span. To cope this linear growth of the variance, we use the regular $n^{-1/2}$ -normalization and the realized martingale differences from (5.7) to construct the test statistic

$$\hat{T} \equiv \max_{j \in [p]} \frac{1}{\sqrt{n}} \sum_{t=1}^n \hat{d}_{tj}. \quad (5.8)$$

Unfortunately, the correlation structure of the martingale differences across decision makers is unknown, and hence the computation of the critical value is not straightforward. In order to overcome this obstacle, we use block multiplier bootstrap method, with the big-block-small-block idea that dates back to Bernstein (1927). We note that this block scheme has been used in Yu (1994); Chen et al. (2016); Chernozhukov et al. (2019) to prove various convergence results under different weak dependencies.

The big-block-small-block idea is simple. It splits the entire sequence of the martingale differences into big blocks and small blocks alternatively, and drops the small blocks away. When the interdependency of the original sequence decays fast enough in terms of the distance between blocks, the resulted big blocks can be well approximated by a corresponding sequence of independent random variables. Mathematically, we first choose two block sizes $m_1(n)$ and $m_2(n)$ for big blocks and small blocks respectively, such that as $n \rightarrow \infty$, they satisfy that $m_1(n) = o(n)$ and $m_2(n) = o(m_1(n))$ while $m_2(n) \rightarrow \infty$. Then we set $q = \lfloor n/(m_1 + m_2) \rfloor$ and divide the horizon $\{1, \dots, n\}$ into alternating groups of size m_1 and m_2 . That is, for all $k \in [q]$, we have

$$\begin{aligned} I_k &= \{(k-1)(m_1 + m_2) + 1, \dots, (k-1)(m_1 + m_2) + m_1\}, \\ J_k &= \{(k-1)(m_1 + m_2) + m_1, \dots, k(m_1 + m_2)\}, \end{aligned}$$

and $J_{q+1} = \{q(m_1 + m_2) + 1, \dots, n\}$. Recalling the realized martingale differences from (5.7), we define shorthand $\hat{\mu}_j$ for their sample average over big blocks as

$$\hat{\mu}_j \equiv \sum_{k=1}^q \sum_{t \in I_k} \hat{d}_{tj}.$$

Next, we generate a sequence of i.i.d. standard normal random variables $\{\epsilon_k\}_{k \in [q]}$ that are independent of everything else, and use them together with (5.7) to construct multiplier bootstrap statistic

$$\widehat{W} \equiv \max_{j \in [p]} \left| \frac{1}{\sqrt{m_1 q}} \sum_{k=1}^q \epsilon_k \sum_{t \in I_k} (\hat{d}_{tj} - \hat{\mu}_j) \right|.$$

Then for any given level of significance α , we use the critical value

$$\widehat{c}(\alpha) \equiv (1 - \alpha)\text{-quantile of } \widehat{W} \text{ given } \widehat{\mathbf{d}}_1. \quad (5.9)$$

That is, we reject the null hypothesis H_0 if our test statistic from (5.8) $\widehat{T} > \widehat{c}(\alpha)$. Based on the critical value $\widehat{c}(\alpha)$, we also have the simultaneous confidence intervals for the total rewards $\mathbf{R}_n \in \mathbb{R}^p$ under H_0 given as

$$I_j(n, \alpha) = \left[v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} - \widehat{c}(\alpha)\sqrt{n}, v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} + \widehat{c}(\alpha)\sqrt{n} \right] \quad \text{for all } j \in [p].$$

5.3.2 Asymptotic Validity

The asymptotic validity of the block multiplier bootstrap procedure is given in the following theorem, which is adopted from Chernozhukov et al. (2019, Theorem B.1).

Theorem 5.1 (Testing for MDPs with regenerative states). Consider p (possibly) correlated MDPs described as in Section 5.2. Suppose that there exist constants c_0, c_1, c_2, C_1 such that the following conditions hold.

(i) *Uniform bound on the martingale differences.*

$$\|\mathbf{d}_t\|_\infty \leq C_1 < \infty \quad \text{and} \quad \min_{j \in [p]} \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] \geq c_1, \quad \text{for all } t \in [n].$$

(ii) *Tail of regeneration time. The MDPs satisfy the regenerative property defined in Definition 5.1, with tail of the regeneration time satisfying*

$$\mathbb{P}\{\tau > k\} \leq \varphi(p, k) \leq p \exp(-c_0 k).$$

(iii) *The parameters satisfy:*

$$\max \left\{ \frac{np}{m_1 + m_2} \exp(-c_0 m_2), \frac{m_2}{m_1} \log^2 p \right\} \leq C_1 n^{-c_2}, \quad \text{and} \quad m_1 \log^{5/2}(pn) \leq C_1 n^{1/2-c_2}. \quad (5.10)$$

Then under the null hypothesis H_0 in (5.1), we have that the critical value (5.9) generated by the block multiplier bootstrap satisfies that

$$\left| \mathbb{P}\{\hat{T} > \hat{c}(\alpha)\} - \alpha \right| \leq C n^{-c},$$

with the two constants c, C depending only on c_0, c_1, c_2, C_1 .

5.4 MDPs with Absorbing States

The other special class of MDPs is that with *absorbing state*. Such absorbing state usually makes the variance of the total reward collected by the optimal policy and near-optimal policies grow sublinearly in the length of time horizon. Specifically, recall the martingale differences given in (5.5), that is,

$$\sum_{t=1}^n \mathbf{d}_t = \mathbf{R}_n - \mathbb{E}[\mathbf{R}_n].$$

In this section we focus on MDPs in which the variance of R_{nj} grows sublinearly in n for all $j \in [p]$. The subtle differences between this class of MDPs and the class of MDPs with regenerative property make our previous analysis inapplicable. Fortunately, however, when more knowledge of the martingale structure is available, one can still conduct bootstrap procedure to obtain critical value for simultaneous hypothesis testings and to construct simultaneous confidence intervals for the total rewards collected by each decision maker that is implementing her designated policy.

5.4.1 Test Statistics and the Multiplier Bootstrap

We start with constructing test statistics with *self-normalization*. Let $\sigma_{nj}^2 = n^{-1} \text{Var} \{ \sum_{t=1}^n d_{tj} \}$ be the normalized variance. Then we normalize the martingale differences by square root of this variance to obtain

$$Z_{tj} \equiv \frac{d_{tj}}{\sigma_{nj}}, \quad \text{for all } t \in [n] \text{ and } j \in [p],$$

so that the associated variance is

$$\text{Var} \left\{ \sum_{t=1}^n Z_{tj} \right\} = \frac{1}{\sigma_{nj}^2} \text{Var} \left\{ \sum_{t=1}^n d_{tj} \right\} = n, \text{ for all } j \in [p].$$

Then from the realized martingale differences (5.7), we obtain the sample variance

$$\hat{\sigma}_{nj}^2 \equiv \frac{1}{n} \sum_{t=1}^n \hat{d}_{tj}^2 \quad \text{for all } j \in [p].$$

With this last sample variance, we can further obtain the realized self-normalized quantity

$$\hat{Z}_{tj} \equiv \frac{\hat{d}_{tj}}{\hat{\sigma}_{nj}}, \quad \text{for all } t \in [n] \text{ and } j \in [p].$$

In principle, with the knowledge of all marginal distributions of the exogenous uncertainty of the possibly correlated MDPs, it is possible to know the exact marginal distributions of the vector of total rewards $\mathbf{R}_n \in \mathbb{R}^p$. However, the covariance structure of \mathbf{R}_n is unknown since the correlation structure of the exogenous uncertainty is unknown. To overcome this, we use the following bootstrap procedure to generate the critical value. We generate i.i.d. standard normal random variables $\{\epsilon_t\}_{t \in [n]}$ that are independent from everything else, and let

$$\hat{Z} = \max_{j \in [p]} \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \epsilon_t \hat{Z}_{tj} \right|.$$

Then we set the critical value as

$$\hat{c}(\alpha) \equiv (1 - \alpha)\text{-quantile of } \hat{Z} \text{ conditional on } \{\hat{Z}_{1j}\}_{j \in [p]}.$$

Based on this test statistics, we reject the null hypothesis H_0 if our test statistics

$$\hat{T} \equiv \max_{j \in [p]} \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \hat{Z}_{tj} \right| > \hat{c}(\alpha).$$

Moreover, we have the simultaneous confidence intervals for the vector of total rewards $\mathbf{R}_n \in \mathbb{R}^p$, given as

$$\left[v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} - \hat{c}(\alpha) \sqrt{n} \hat{\sigma}_{nj}, v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} + \hat{c}(\alpha) \sqrt{n} \hat{\sigma}_{nj} \right] \quad \text{for all } j \in [p]. \quad (5.11)$$

5.4.2 Asymptotic Validity

The above bootstrap procedure for hypothesis testing and construction of simultaneous confidence intervals is asymptotically valid under suitable conditions. The asymptotic validity is mainly based on Belloni and Oliveira (2018, Assumption 2), and we tailor the conditions to adapt to our setup and notations. As we will see in Section 5.6, it is usually not trivial to verify these conditions in applications, and such verifications often require a substantial understanding of the MDPs at hand. The asymptotic validity result is made precise in the following statement.

Theorem 5.2 (Testing for MDPs with absorbing states). Let δ_n and ψ_n be two fixed sequences that are going to zero as n increases. Suppose that the martingale differences defined in (5.6) satisfy the following conditions.

- (i) $\frac{1}{n} \sum_{t=1}^n \mathbb{E} \left[\max_{j \in [p]} |Z_{tj}|^3 \right] \leq n^\rho$ for some $\rho \in (0, 1/2)$;
- (ii) The matrices $V \equiv \frac{1}{n} \sum_{t=1}^n \mathbb{E} [\mathbf{Z}_t \mathbf{Z}_t']$ and $V_n \equiv \frac{1}{n} \sum_{t=1}^n \mathbb{E} [\mathbf{Z}_t \mathbf{Z}_t' \mid \mathcal{F}_{t-1}]$ satisfy that

$$\max_{i,j \in [p]} |V_{ij}| \leq C \quad \text{and} \quad \mathbb{E} \left[\max_{i,j \in [p]} |V_{n,ij} - V_{ij}| \right] \leq \frac{\delta_n}{\log^2(np)},$$

and with probability at least $1 - \psi_n$, we have that

$$V_n \leq (1 + \delta_n / \log^2 p) V, \quad \text{and} \quad \max_{i,j \in [p]} |V_{n,ij} - V_{ij}| \leq \frac{\delta_n}{\log^2(np)};$$

(iii) With probability at least $1 - \psi_n$, we have that

$$\max_{i,j \in [p]} \left| \frac{1}{n} \sum_{t=1}^n Z_{ti} Z_{tj} - \mathbb{E}[Z_{ti} Z_{tj} \mid \mathcal{F}_{t-1}] \right| \vee \max_{j \in [p]} \left| \frac{1}{n} \sum_{t=1}^n (\hat{Z}_{tj} - Z_{tj})^2 \right| \leq \frac{\delta_n^2}{\log^2(np)};$$

(iv) $\psi_n \leq \delta_n^{1/2}$, $\log^{7/2} p \leq \delta_n n^{1/2-\rho}$, and $\min_{j \in [p]} \sigma_j \geq c$.

Then the simultaneous confidence intervals in (5.11) are asymptotic valid. That is, as $n \rightarrow \infty$,

$$\left| \mathbb{P} \left\{ \sum_{t=1}^n \hat{d}_{tj} - \hat{c}(\alpha) \sqrt{n} \hat{\sigma}_{nj} \leq R_{nj} - v_{nj}(x_{0j}) \leq \sum_{t=1}^n \hat{d}_{tj} + \hat{c}(\alpha) \sqrt{n} \hat{\sigma}_{nj}, \text{ for all } j \in [p] \right\} - (1 - \alpha) \right| \rightarrow 0.$$

5.5 Application: Testing Stationary Inventory Control Policies

To illustrate how our MDP framework and the multiplier bootstrap procedure in Section 5.3 can be applied to concrete applications, we consider *dynamic inventory control problems*. The setup we will be focusing on is in which the decision makers are implementing the stationary policies that are optimal to the infinite horizon problem. While we observe a time horizon of finite length, and we are interested in simultaneously testing whether each of the decision makers is implementing her designated policy.

5.5.1 Problem Setup (Scarf, 1959; Iglehart, 1963).

We now describe a dynamic inventory control problem faced in parallel by multiple decision makers. Suppose there are p decision makers each of whom manages a inventory of a single product over an infinite time horizon, and we get to observe the costs they incur and their each-day inventory levels over a time horizon of the length n . For each decision maker $j \in [p]$, the inventory level starts from initial level X_{0j} , and she implements a stationary policy that does not change over time. Then in each time period $t \in [n]$, she needs to decide whether to order the product or not, and how many units of the product to order if she decides to place an order. If we denote $\gamma_j(x)$ as her order-up-to function, that is, she will order $(\gamma_j(x) - x)$ units of the product at time period t if her inventory level from the last period is $X_{(t-1)j} = x$. Then it is assumed that the order will be fulfilled right away and an uncertain demand ξ_{tj} will be realized after the order is delivered (if she places an order), so

that the inventory level evolves as

$$X_{tj} = \gamma_j(X_{(t-1)j}) - \xi_{tj}. \quad (5.12)$$

The marginal distributions of ξ_{tj} 's for all $j \in [p]$ are known. However, the correlation among different coordinates of the vector $\boldsymbol{\xi}_t$ of demands is unknown and could be arbitrary. There are three types of costs. The ordering cost

$$c_{o,j}(x) = \begin{cases} K + c_j x & \text{if } x > 0 \\ 0 & \text{if } x = 0, \end{cases} \quad (5.13)$$

the per-unit holding cost $c_{h,j}$, and the per-unit backlog penalty $c_{p,j}$. More specifically, there is a fixed ordering cost K , and the j -th individual decision maker needs to pay $K + c_j(y - x)$ to bring the inventory of her product from x up to y before the demand realization. Then after the demand is realized, if the remaining inventory level x is positive, then she incurs a holding cost $c_{h,j}x$ and these x units of product will be carried over to the next period; while if the inventory level x is negative, she incurs a backlog penalty $-c_{p,j}x$ for not being able to meet these demands. The latter two types of costs can be written in one single function as

$$L_j(x) = \begin{cases} c_{h,j}x & \text{if } x \geq 0 \\ -c_{p,j}x & \text{if } x < 0. \end{cases} \quad (5.14)$$

All the decision makers are implementing the optimal stationary base-stock policy for the infinite horizon problem (Iglehart, 1963). That is, for each decision maker $j \in [p]$, there is a pair of base-stock levels (s_j, S_j) with $s_j < S_j$ such that whenever the inventory falls below s_j , she will place an order and bring the inventory level up to S_j . Equivalently, she has the order-up-to function given by

$$\gamma_j(x) = \begin{cases} S_j & \text{if } x \leq s_j \\ x & \text{if } x > s_j. \end{cases} \quad (5.15)$$

The total cost she incurs by implementing this base-stock policy over a time horizon of length n is given as

$$R_{nj} = \sum_{t=1}^n \{c_{o,j}(\gamma_j(X_{(t-1)j}) - X_{(t-1)j}) + L_j(X_{tj})\}.$$

We also have the value function representation

$$v_{nj}(X_{0j}) = \mathbb{E}[R_{nj}],$$

and for all $t \in [n]$,

$$v_{(n-t)j}(X_{t,j}) = \mathbb{E} \left[\sum_{\ell=t}^n \{c_{o,j}(\gamma_j(X_{(\ell-1)j}) - X_{(\ell-1)j}) + L_j(X_{\ell j})\} \right].$$

As a result, our martingale differences in this context are

$$\begin{aligned} d_{tj} &= \mathbb{E}[R_{nj} \mid \mathcal{F}_t] - \mathbb{E}[R_{nj} \mid \mathcal{F}_{t-1}] \\ &= c_{o,j}(\gamma_j(X_{(t-1)j}) - X_{(t-1)j}) + L(X_{tj}) + v_{(n-t)j}(X_{tj}) - v_{(n-t+1)j}(X_{(t-1)j}). \end{aligned}$$

With the shorthand $\hat{v}_k(x) = L(x) + v_k(x)$, one can easily obtain that

$$d_{tj} = \hat{v}_{(n-t)j}(\gamma_j(X_{(t-1)j}) - \xi_{tj}) - \mathbb{E}[\hat{v}_{(n-t)j}(\gamma_j(X_{(t-1)j}) - \xi_{tj}) \mid \mathcal{F}_{t-1}]. \quad (5.16)$$

The marginal distributions of demands play a crucial role in the analysis, and we put the certain regularity conditions on these marginal distributions. For all the base-stock level pairs (s_j, S_j) , $j \in [p]$, let $\underline{s} = \min_{j \in [p]} s_j$ and $\bar{s} = \max_{j \in [p]} S_j - s_j$. The regularity conditions we need are the following.

Definition 5.2 (Regular class of distributions). A p -dimensional non-negative bounded random variable $\mathbf{W} \in \mathbb{R}^p$ with continuous distribution $F(\cdot) : [0, \bar{D}]^p \rightarrow [0, 1]$ and marginal densities f_j for all $j \in [p]$, is said to be *regular* if the following conditions are satisfied.

- Finite, non-vanishing means:

$$\underline{\mu} \equiv \min_{j \in [p]} \mathbb{E}[W_j], \quad \bar{\mu} \equiv \max_{j \in [p]} \mathbb{E}[W_j], \quad \text{and} \quad 0 < \underline{\mu} \leq \bar{\mu} < \infty.$$

- Finite variance:

$$\max_{j \in [p]} \text{Var}\{W_j\} \leq \bar{\sigma}^2 < \infty.$$

- Non-flat marginal density: for all the marginal distribution F_j with density f_j ,

$$\max_{j \in [p]} \{\mathbb{P}\{W_j \leq S_j\}, \mathbb{P}\{W_j \geq s_j\}\} \leq \bar{p} < 1 \quad \text{and} \quad \min_{j \in [p]} \inf_{w \in [0, \bar{D}]} f_j(w) \geq \eta > 0. \quad (5.17)$$

- For each marginal distribution with density f_j , and for each $\epsilon \geq 0$, there is $\hat{w} = \hat{w}(\epsilon)$ such that

$$f_j(w) - f_j(w + \epsilon) \leq 0 \quad \text{for all } w \leq \hat{w}, \text{ and}$$

$$f_j(w) - f_j(w + \epsilon) \geq 0 \quad \text{for all } w \geq \hat{w}.$$

The boundedness of the means and variances essentially rules out trivial cases in which demands are deterministic. The last condition is borrowed from Arlotto and Steele (2016, Definition 12), which is used to guarantee boundedness of the contraction coefficient there. We will use similar analysis on the individual inventory levels.

If the joint uncertain demands faced by each of the decision makers satisfy the regularity conditions given in Definition 5.2, then we can apply the bootstrap procedure from Section 5.3 to construct simultaneous confidence intervals for the total costs incurred by each decision maker with asymptotic validity. Such results are summarized in the following theorem.

Theorem 5.3. If the joint demand distribution F is regular as defined in Definition 5.2, then the simultaneous confidence intervals

$$\prod_{j=1}^p \left[v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} - \hat{c}(\alpha)\sqrt{n}, v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} + \hat{c}(\alpha)\sqrt{n} \right]$$

as constructed in Section 5.3 are asymptotic valid. That is,

$$\mathbb{P} \left\{ v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} - \hat{c}(\alpha)\sqrt{n} \leq R_{nj} \leq v_{nj}(x_{0j}) + \sum_{t=1}^n \hat{d}_{tj} + \hat{c}(\alpha)\sqrt{n} \right\} \rightarrow (1 - \alpha), \quad \text{as } n \rightarrow \infty.$$

5.5.2 Regeneration Time with Exponential Tail

We consider the regeneration time $(\tau + 1)$ which is the first time that every decision maker has at least put an order. We establish an exponentially decaying bound for the tail probability of τ , which is a stopping time of the Markov chain $\{\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_n\} \subset \mathbb{R}^p$ of the aggregated inventory levels induced by each individual decision maker's policy. The basic idea is that each time a decision maker places an order, her inventory level regenerates in the sense that the future evolution of her inventory level becomes independent of the past history. If one looks at the particular stopping time τ , which is the first time that for every decision maker j , her inventory level has fall below s_j for at least once, then some moments conditions on the marginal distributions of the uncertainty imply an exponentially decaying right tail probability of τ .

To establish this exponential decaying tail, we need the following lemma which concerns lower tail probability of sum of independent non-negative random variables. It can be obtained from several different well-known concentration inequalities. We provide a proof for the sake of completeness.

Lemma 5.1 (Lower tail of independent non-negative random variables). *Let X_1, \dots, X_n be independent non-negative random variables. Then for any $t > 0$,*

$$\mathbb{P} \left\{ \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \leq -t \right\} \leq \exp \left(\frac{-t^2}{2 \sum_{i=1}^n \mathbb{E}[X_i^2]} \right).$$

Proof. For any non-negative random variable X , the simple inequality $e^{-u} \leq 1 - u + u^2/2$ for all $u \geq 0$ tells us that for all $\lambda > 0$,

$$\mathbb{E}[e^{-\lambda X}] \leq 1 - \lambda \mathbb{E}[X] + \lambda^2 \mathbb{E}[X^2]/2.$$

Another elementary inequality $\log u \leq u - 1$ for all $u > 0$ gives us that

$$\log \mathbb{E}[e^{-\lambda X}] \leq \mathbb{E}[e^{-\lambda X}] - 1,$$

which in turn implies that

$$\log \mathbb{E}[e^{-\lambda X}] \leq -\lambda \mathbb{E}[X] + \lambda^2 \mathbb{E}[X^2]/2.$$

If we take exponential of both sides and rearrange, we obtain that for all $\lambda > 0$,

$$\mathbb{E}[e^{\lambda(X - \mathbb{E}[X])}] \leq e^{\lambda^2 \mathbb{E}[X^2]/2}.$$

This last inequality tells us that when we have independent non-negative random variables X_1, \dots, X_n , the independence implies that

$$\mathbb{E}[e^{\lambda \sum_{i=1}^n (X_i - \mathbb{E}[X_i])}] = \prod_{i=1}^n \mathbb{E}[e^{\lambda(X_i - \mathbb{E}[X_i])}] \leq e^{\lambda^2 \sum_{i=1}^n \mathbb{E}[X_i^2]/2}.$$

By Markov inequality, we know that for all $t > 0$,

$$\mathbb{P} \left\{ \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \leq -t \right\} \leq e^{-\lambda t} \mathbb{E}[e^{\lambda \sum_{i=1}^n (X_i - \mathbb{E}[X_i])}] \leq \exp \left\{ -\lambda t + \lambda^2 \sum_{i=1}^n \mathbb{E}[X_i^2]/2 \right\}.$$

Note that the last right-hand side is minimized by choosing $\lambda = t / \sum_{i=1}^n \mathbb{E}[X_i^2]$, which in turn gives us the final tail probability bound

$$\mathbb{P} \left\{ \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \leq -t \right\} \leq \exp \left(\frac{-t^2}{2 \sum_{i=1}^n \mathbb{E}[X_i^2]} \right),$$

just as we needed. □

With this last lemma, we are able to obtain an exponential-decay upper bound on the right tail probability of the stopping time τ . In fact, it is without loss of generality to assume the initial inventory $X_{0j} = S_j$ for all $j \in [p]$. General values of X_{0j} can be analyzed similarly. The tail bound is given in the following proposition.

Proposition 5.1 (Tail bound of the mixing time of joint inventory levels). *There exist constants K_0 and \bar{s} that depend only on the marginal distributions of $F(\cdot) : \mathbb{R}_+^p \rightarrow [0, 1]$ such that for $k \geq K_0$,*

$$\mathbb{P}\{\tau > k\} \leq p \exp\left(-\frac{(k\underline{\mu} - \bar{s})^2}{2k(\bar{\mu}^2 + \bar{\sigma}^2)}\right),$$

with the constants $\underline{\mu}, \bar{\mu}, \bar{\sigma}^2$ given in Definition 5.2.

Proof. Recall that the structure of the (s, S) base-stock policy, and that $\bar{s} \equiv \max_{j \in [p]}(S_j - s_j)$, we have the inclusion

$$\{\omega : \tau > k\} \subset \bigcup_{j \in [p]} \left\{ \omega : \sum_{t=1}^k \xi_{tj} < \bar{s} \right\}. \quad (5.18)$$

Therefore, we have the tail probability bound

$$\mathbb{P}\{\tau > k\} \leq \mathbb{P}\left\{ \bigcup_{j \in [p]} \left\{ \sum_{t=1}^k \xi_{tj} < \bar{s} \right\} \right\} \leq \sum_{j \in [p]} \mathbb{P}\left\{ \sum_{t=1}^k \xi_{tj} < \bar{s} \right\}.$$

For each summand, we apply Lemma 5.1 to obtain that

$$\mathbb{P}\left\{ \sum_{t=1}^k \xi_{tj} < \bar{s} \right\} = \mathbb{P}\left\{ \sum_{t=1}^k (\xi_{tj} - \mu_j) < -(k\mu_j - \bar{s}) \right\} \leq \exp\left(-\frac{(k\mu_j - \bar{s})^2}{2k(\mu_j^2 + \sigma_j^2)}\right),$$

where in the last step a lower-bound on $k \geq K_0$ with $K_0 \equiv \bar{s}/\underline{\mu}$ is needed ($\underline{\mu}$ defined in Definition 5.2) to guarantee the positiveness $k\mu_j - \bar{s} > 0$. In other words, this tail probability bound is valid for $k \geq K_0$. Finally, some simple algebra gives us the inequality

$$-\frac{(k\mu_j - \bar{s})^2}{2k(\mu_j^2 + \sigma_j^2)} \leq -\frac{(k\underline{\mu} - \bar{s})^2}{2k(\bar{\mu}^2 + \bar{\sigma}^2)},$$

which completes the proof. \square

Now we are left with proving that upon regeneration, \mathbf{X}_τ becomes independent of the initial state \mathbf{X}_0 . Note that under the base-stock policy, the states $\{\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_n\} \subset \mathbb{R}^p$ of the inventory levels form a Markov chain that evolves as that for all $t \in [n]$ and all $j \in [p]$,

$$X_{tj} = \begin{cases} X_{(t-1)j} - \xi_{tj} & \text{if } X_{(t-1)j} > s_j \\ S_j - \xi_{tj} & \text{if } X_{(t-1)j} \leq s_j. \end{cases} \quad (5.19)$$

In addition to the stopping time τ , the base-stock policy also induces, for each decision maker $j \in [p]$, a sequence of stopping times $0 \equiv \kappa_0^j < \kappa_1^j < \dots < \kappa_k^j \leq n$, at which orders are placed. Specifically, by setting $\kappa_0^j = 0$, one can recursively define the stopping times

$$\kappa_k^j = \max \left\{ t : \kappa_{k-1}^j < t \leq n \text{ and } \sum_{\ell=\kappa_{k-1}^j}^t \xi_{\ell j} \leq S_j - s_j \right\} \quad \text{for } k \geq 1. \quad (5.20)$$

In addition to these stopping times, we also have the equivalent representation of τ as

$$\tau = \min \left\{ t \in [n] \cup \{\infty\} : \kappa_1^j \geq t \text{ for all } j \in [p] \right\}. \quad (5.21)$$

Then for each $m \in [n]$, we set $\boldsymbol{\eta}(m) \in \mathbb{R}^p$ as

$$\eta_j(m) = \max \left\{ \kappa_k^j \leq m : k \geq 0 \right\} \quad \text{for all } j \in [p]. \quad (5.22)$$

In words, the j -th element of the vector $\boldsymbol{\eta}(m)$ is the last time before m that the inventory level of the j -th decision maker falls below s_j . For a fixed integer $m \geq 1$, we define the configuration set $\Gamma(m)$ as

$$\Gamma(m) \equiv \left\{ \mathbf{t} \in [m]^p : \max_{j \in [p]} t_j = m \right\}.$$

That is, the elements of $\Gamma(m)$ all have the form $\mathbf{t} = (t_j)_{j \in [p]}$ where $t_j \in [m]$ for all $j \in [p]$ and the largest coordinate of \mathbf{t} is equal to m . We call an element $\mathbf{t} \in \Gamma(m)$ a configuration. The structure of the base-stock policy implies that, with the knowledge of $\{\tau = m, \boldsymbol{\eta}(m) = \mathbf{t}\}$, the state \mathbf{X}_m becomes independent of the initial state \mathbf{X}_0 . This property is summarized in the following lemma.

Lemma 5.2. For all $m \in [n-1]$, all configuration $\mathbf{t} \in \Gamma(m)$, and all $\mathbf{x}_0 \in \mathbb{R}^p$, if $\boldsymbol{\eta}(m) \in [m]^p$ is the vector of the last time each decision maker puts an order, then

$$\mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t} \mid \mathbf{X}_0 = \mathbf{x}_0\} = \mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t}\}, \quad \text{for all } B \subset \mathbb{R}^p.$$

Proof. The structure of the base-stock policy (5.19) tells us that, for each $j \in [p]$, whenever there is an order placed at some stopping time κ as defined in (5.20) (that is, $X_{\kappa j} \leq s_j$), the next-period inventory level depends only on the base-stock level S_j and the demand $\xi_{(\kappa+1)j}$. Moreover, the state X_{tj} before the next backorder is independent of X_{0j} in the sense that for all $k \geq 1$, all $\kappa_{k-1}^j < t \leq \kappa_k^j$, and all $B \subset (s_j, \infty)$,

$$\mathbb{P}\{X_{tj} \in B, \kappa_{k-1}^j = \kappa \mid X_{0j}\} = \mathbb{P}\left\{\sum_{\ell=\kappa+1}^t \xi_{\ell j} \in S_j - B\right\} = \mathbb{P}\{X_{tj} \in B, \kappa_{k-1}^j = \kappa\}.$$

Similarly, for all $m \in [n-1]$, and any $\mathbf{t} \in \Gamma(m)$, we have that for all $B \subset \mathbb{R}^p$,

$$\begin{aligned} \mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t} \mid \mathbf{X}_0\} &= \mathbb{P}\left\{\sum_{\ell=t_j}^m \xi_{\ell j} \in S_j - B_j, j \in [p]\right\} \\ &= \mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t}\}, \end{aligned}$$

so the proof is complete. \square

With this last lemma, we are able to prove the regenerative property of \mathbf{X}_t at the stopping time τ in (5.21).

Lemma 5.3 (Independence upon regeneration). *The inventory level after the stopping time τ is independent of the initial level \mathbf{X}_0 in the sense that*

$$\mathbb{P}\{\mathbf{X}_{\tau+1} \in B \mid \mathbf{X}_0\} = \mathbb{P}\{\mathbf{X}_{\tau+1} \in B\} \quad \forall B \subset \mathbb{R}^p.$$

Proof. We start with the conditional probability $\mathbb{P}\{\mathbf{X}_{\tau+1} \in B \mid \mathbf{X}_0 = \mathbf{x}_0\}$ and note that by conditioning on τ and $\boldsymbol{\eta}(\tau)$, we have that for all $\mathbf{x}_0 \in \mathbb{R}^p$,

$$\begin{aligned} \mathbb{P}\{\mathbf{X}_{\tau+1} \in B \mid \mathbf{X}_0 = \mathbf{x}_0\} &= \sum_{m=1}^{\infty} \sum_{\mathbf{t} \in \Gamma(m)} \mathbb{P}\{\mathbf{X}_{\tau+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t} \mid \mathbf{X}_0 = \mathbf{x}_0\} \\ &= \sum_{m=1}^{\infty} \sum_{\mathbf{t} \in \Gamma(m)} \mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t} \mid \mathbf{X}_0 = \mathbf{x}_0\}. \end{aligned}$$

From Lemma 5.2 we know that each summand in this last expression satisfies that

$$\mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t} \mid \mathbf{X}_0 = \mathbf{x}_0\} = \mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t}\}.$$

Combining the last two identities, we obtain that

$$\mathbb{P}\{\mathbf{X}_{\tau+1} \in B \mid \mathbf{X}_0 = \mathbf{x}_0\} = \sum_{m=1}^{\infty} \sum_{\mathbf{t} \in \Gamma(m)} \mathbb{P}\{\mathbf{X}_{m+1} \in B, \tau = m, \boldsymbol{\eta}(m) = \mathbf{t}\} = \mathbb{P}\{\mathbf{X}_{\tau+1} \in B\},$$

as we wanted. \square

5.5.3 Analysis on the Martingale Differences

In addition to the regenerative property of the inventory level, we also need to verify other conditions in Theorem 5.1 to prove Theorem 5.3. The first one is regarding the boundedness of the martingale differences.

Proposition 5.2 (Martingale difference upper bound). *There exists a constant B such that the martingale differences are uniformly upper-bounded*

$$\|d_{t,j}\|_{\infty} \leq B, \quad \text{for all } j \in [p] \text{ and all } t \in [n].$$

The proof of this last boundedness property relies on the following intermediate step. This intermediate step essentially draws from Arlotto and Steele (2016, Lemma 14 and Lemma 15).

Proposition 5.3. *When the demand distribution satisfies the conditions in Definition 5.2, we have that for all $j \in [p]$ and all $t \in [n]$,*

$$\sup_{B \subset \mathbb{R}} \left| \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x\} - \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x'\} \right| \leq \bar{p} < 1,$$

with the constant \bar{p} from Definition 5.2.

Proof. If we define $B_x \equiv \gamma_j(x) - B$, then from (5.12) and (5.15) we know that

$$\mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x\} = \mathbb{P}\{\xi_{1j} \in B_x\} = \int_{B_x} f_j(w) dw.$$

As a result, if we assume $x \leq x'$ and set $\epsilon = \gamma_j(x') - \gamma_j(x) \geq 0$, we obtain that

$$\mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x\} = \mathbb{P}\{\xi_{1j} - \epsilon \in B_x\} = \int_{B_x} f_j(w + \epsilon) dw.$$

Combining the last two identities, we further obtain that

$$\begin{aligned} & \sup_{B \subset \mathbb{R}} \left| \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x\} - \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x'\} \right| \\ &= \sup_{B \subset \mathbb{R}} \left| \int_{B_x} f_j(w) dw - \int_{B_x} f_j(w + \epsilon) dw \right| = \left| \int_{B_x^*} f_j(w) - f_j(w + \epsilon) dw \right|, \end{aligned}$$

where $B_x^* \equiv \{w : f_j(w) \geq f_j(w + \epsilon)\}$. Therefore, from Definition 5.2 we know that $B_x^* = [\hat{w}, \bar{D}]$, so

$$\begin{aligned} & \sup_{B \subset \mathbb{R}} \left| \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x\} - \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x'\} \right| \\ &= \int_{\hat{w}}^{\bar{D}} f_j(w) - f_j(w + \epsilon) dw = \mathbb{P}\{\hat{w} \leq \xi_{1j} \leq \hat{w} + \epsilon\}. \end{aligned}$$

Since $\epsilon = \gamma_j(x') - \gamma_j(x) \leq S_j - s_j$, if $\hat{w} + \epsilon \leq S_j$, then $\mathbb{P}\{\hat{w} \leq \xi_{1j} \leq \hat{w} + \epsilon\} \leq \mathbb{P}\{\xi_{1j} \leq S_j\}$, which is bounded by \bar{p} by Definition 5.2. If $\hat{w} + \epsilon > S_j$, then $\hat{w} > s_j$, and so $\mathbb{P}\{\hat{w} \leq \xi_{1j} \leq \hat{w} + \epsilon\} \leq \mathbb{P}\{\xi_{1j} > s_j\}$, which is again bounded by \bar{p} by Definition 5.2. Therefore, we obtain that

$$\sup_{B \subset \mathbb{R}} \left| \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x\} - \mathbb{P}\{X_{(t+1)j} \in B \mid X_{tj} = x'\} \right| \leq \bar{p} < 1,$$

completing the proof. \square

Now we are ready to prove Proposition 5.2.

Proof. (Proof of Proposition 5.2.) We first verify that the one-period cost of the dynamic inventory control problem is uniformly bounded. In fact, for the j -th decision maker, the cost she incurs during period t is given as

$$\mathcal{P}_{tj} = c_{o,j}(\gamma_j(X_{(t-1)j}) - X_{(t-1)j}) + L_j(\gamma_j(X_{(t-1)j}) - \xi_{tj}).$$

If we recall the ordering cost (5.13), the holding-and-shortage cost (5.14), and the order-up-to function (5.15), we obtain that

$$\mathcal{P}_{tj} = \begin{cases} c_{h,j}(X_{(t-1)j} - \xi_{tj}) & \text{if } s_j < X_{(t-1)j} \text{ and } \xi_{tj} \leq X_{(t-1)j} \\ c_{p,j}(\xi_{tj} - X_{(t-1)j}) & \text{if } s_j < X_{(t-1)j} < \xi_{tj} \\ K + c_{o,j}(S_j - X_{(t-1)j}) + c_{h,j}(S_j - \xi_{tj}) & \text{if } X_{(t-1)j} \leq s_j \text{ and } \xi_{tj} \leq S_j \\ K + c_{o,j}(S_j - X_{(t-1)j}) + c_{p,j}(\xi_{tj} - S_j) & \text{if } X_{(t-1)j} \leq s_j \text{ and } S_j < \xi_{tj}. \end{cases}$$

Since the demand ξ_{tj} is uniformly bounded within the interval $[0, \bar{D}]$, we can easily verify that the one-period cost satisfies that

$$0 \leq \mathcal{P}_{tj} \leq \max\{c_{h,j}S_j, c_{p,j}(\bar{D} - S_j), K + c_{o,j}(S_j - s_j) + c_{h,j}S_j, K + c_{o,j}(S_j - s_j) + c_{p,j}(\bar{D} - S_j)\} \equiv \bar{C}_j.$$

Then we set $\bar{C} = \max_{j \in [p]} \bar{C}_j$, and apply Arlotto and Steele (2016, Lemma 8) and Proposition 5.3 to obtain that the martingale differences satisfy that

$$\|d_{tj}\|_{\infty} \leq \frac{M\bar{C}}{\bar{p}} \equiv B \quad \text{for all } t \in [n] \text{ and all } j \in [p],$$

completing the proof. \square

Next, we work towards an lower bound on the conditional second moment of martingale differences, $\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]$. This lower bound is summarized in the following lemma, whose proof relies on a coupling argument.

Lemma 5.4 (Martingale difference second moment lower bound). There exists a constant $\ell_1 > 0$ such that the conditional second moment of the martingale differences are uniformly bounded from below, that is,

$$\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] \geq \ell_1, \quad \text{for all } j \in [p] \text{ and all } t \in [n].$$

As a result, for any set $I = \{i+1, i+2, \dots, i+q\}$, one has the following variance lower bound.

$$\text{Var} \left\{ \sum_{t \in I} d_{tj} \right\} \geq \ell_1 q.$$

Proof. Let $\{\tilde{\xi}_{1j}, \dots, \tilde{\xi}_{nj}\}$ be an independent copy of the (actual) demands $\{\xi_{1j}, \dots, \xi_{nj}\}$ and we use coupling to bound the conditional second moment $\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]$. Note that the inventory level $X_{(t-1)j}$ at the beginning of time period t is \mathcal{F}_{t-1} -measurable, i.e., it is totally determined by the demand realizations $\{\xi_{ij}\}_{i \in [t-1]}$. While the next period inventory level X_{tj} depends on $X_{(t-1)j}$ as well as the demand realization ξ_{tj} through the relation $X_{tj} = \gamma_j(X_{(t-1)j}) - \xi_{tj}$. Recall that for any two i.i.d. zero-mean random variables $Y \sim Z$, we have that

$$\text{Var}\{Y\} = \mathbb{E}[Y^2] = \frac{1}{2} \mathbb{E}[(Y - Z)^2].$$

If we apply this to our martingale difference d_{tj} from (5.16), we would get

$$\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] = \frac{1}{2} \mathbb{E} \left[\left\{ \hat{v}_{(n-t)j}(\gamma_j(X_{(t-1)j}) - \xi_{tj}) - \hat{v}_{(n-t)j}(\gamma_j(X_{(t-1)j}) - \tilde{\xi}_{tj}) \right\}^2 \mid \mathcal{F}_{t-1} \right].$$

Next, we restrict ourselves on the event

$$G = G(X_{(t-1)j}) \equiv \left\{ \omega : \xi_{tj}, \tilde{\xi}_{tj} \in [\gamma_j(X_{(t-1)j}) - s_j, \gamma_j(X_{(t-1)j})] \right\}.$$

In fact, the choice of event G would guarantee that a new order will be placed for both of the two trajectories associated with $\{\xi_{tj}\}$ and $\{\tilde{\xi}_{tj}\}$. To see this, one only needs to verify that the inventory level X_{tj} at the end of time period t is below the threshold s_j , that is,

$$\gamma_j(X_{(t-1)j}) - \xi_{tj} = X_{tj} \in [0, s_j].$$

The same holds for the other trajectory with $\{\tilde{\xi}_{tj}\}$. As a result, conditioning on event G , the two trajectories coincide at time period $(t+1)$, and hence all the future randomnesses also become the same. This in turn gives us the lower bound

$$\begin{aligned}\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] &\geq \frac{1}{2} \mathbb{E} \left[(c_j + c_{h,j})^2 \left(\xi_{tj} - \tilde{\xi}_{tj} \right)^2 \mathbb{1}\{G\} \mid \mathcal{F}_{t-1} \right] \\ &= \frac{1}{2} (c_j + c_{h,j})^2 \iint_{\Omega} (y - z)^2 f_j(y) f_j(z) \, dydz,\end{aligned}$$

where the integral region is $\Omega \equiv [\gamma_j(X_{(t-1)j}) - s_j, \gamma_j(X_{(t-1)j})]^2$, and f_j is the marginal density of the demand distribution faced by the j -th individual decision maker. The regularity condition (5.17) of demand distribution implies a strictly positive lower bound on the double integral

$$\iint_{\Omega} (y - z)^2 f_j(y) f_j(z) \, dydz \geq \eta^2 \iint_{\Omega} (y - z)^2 \, dydz = \frac{1}{12} \eta^2 s_j^4,$$

hence

$$\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] \geq \frac{1}{24} (c_j + c_{h,j})^2 \eta^2 s_j^4 \geq \frac{\eta^2}{24} \min_{j \in [p]} \{(c_j + c_{h,j})^2 s_j^4\} \equiv \ell_1.$$

As a result, the variance is lower-bounded by

$$\text{Var} \left\{ \sum_{t \in I} d_{tj} \right\} = \sum_{t=1}^q \mathbb{E} [d_{(t+i)j}^2] = \sum_{t=1}^q \mathbb{E} \left[\mathbb{E} [d_{(t+i)j}^2 \mid \mathcal{F}_{t+i-1}] \right] \geq \ell_1 q,$$

and the proof is complete. \square

With all the above preparations, the proof of Theorem 5.3 becomes straightforward.

Proof. (Proof of Theorem 5.3.) We verify the conditions required in Theorem 5.1.

Condition (i) is verified by Proposition 5.2 and Lemma 5.4.

Condition (ii) is verified by Proposition 5.1 and Lemma 5.3. \square

5.6 Application: Online Monotone Subsequence Selection from a Random Sample

In this section, we consider concrete application of the bootstrap procedure developed in Section 5.4. Specifically, we revisit the sequential monotone subsequence selection problem that was studied in Chapter 3.

5.6.1 Problem Description and Main Result

We remind the readers with the problem faced by a single decision maker described as follows. A sequence of i.i.d. random variables W_1, \dots, W_n from a continuous distribution are sequentially presented to the decision maker, who needs to select a (increasing) subsequence $W_{\tau_1}, \dots, W_{\tau_\ell}$ that satisfies that

$$W_{\tau_1} \leq W_{\tau_2} \leq \dots \leq W_{\tau_\ell},$$

where τ_1, \dots, τ_ℓ are stopping times with respect to the filtration $\mathcal{F}_t = \sigma\{W_1, \dots, W_t\}$, $t \in [n]$. Her objective is to maximize the expected length of the selected subsequence. This problem can be interpreted as a finite-horizon MDP, with the state X_t being the last selected value and the initial state being $X_0 = 0$. When the distribution F of the random variables is continuous, it is without loss of generality to assume that F is the uniform distribution on the unit interval. Moreover, the optimal policy is a time- and state-dependent threshold policy. Specifically, at any time t , let X_{t-1} be the last selected value and let W_t be the observed value. Then there is a threshold function $h_{n-t+1}(\cdot) : [0, 1] \rightarrow [0, 1]$ such that it is optimal to select W_t if and only if

$$X_{t-1} \leq W_t \leq h_{n-t+1}(X_{t-1}).$$

As a result, the state evolves accordingly and can be defined recursively as that for all $t \in [n]$,

$$X_t = \begin{cases} W_t & \text{if } W_t \in [X_{t-1}, h_{n-t+1}(X_{t-1})] \\ X_{t-1} & \text{otherwise.} \end{cases}$$

Then the total reward collected by this optimal policy is given as

$$R_n \equiv \sum_{t=1}^n \mathbb{1}\{W_t \in [X_{t-1}, h_{n-t+1}(X_{t-1})]\},$$

and one also has the optimal value function representation $v_n = \mathbb{E}[R_n]$.

The following upper bound on the length of the selection interval $[X_{t-1}, h_{n-t+1}(X_{t-1})]$ will be useful in our analysis.

Proposition 5.4 (Upper bound on selection interval length). *For any $t \in [n]$, we have that*

$$h_{n-t+1}(X_{t-1}) - X_{t-1} \leq \frac{4 \log(n-t+1) + 6}{\sqrt{2(n-t+1)}}. \quad (5.23)$$

Proof. We apply Theorem 3.1 to obtain a uniform upper bound on the difference $h_{n-t+1}(X_{(t-1)j}) - X_{(t-1)j}$ for all $j \in [p]$. In fact, we know from optimality that for all $j \in [p]$,

$$v_{n-t+1}(X_{(t-1)j}) \leq 1 + v_{n-t+1}(h_{n-t+1}(X_{(t-1)j})) \quad \text{for all } t \in [n].$$

By Theorem 3.1 we know that the left-hand side has lower bound

$$\sqrt{2(n-t+1)(1-X_{(t-1)j})} - 2(\log(n-t+1)+1) \leq v_{n-t+1}(X_{(t-1)j}),$$

and the right-hand side has upper bound

$$1 + v_{n-t+1}(h_{n-t+1}(X_{(t-1)j})) \leq 1 + \sqrt{2(n-t+1)(1-h_{n-t+1}(X_{(t-1)j}))}.$$

Putting these last three inequalities together, we further obtain that

$$\sqrt{2(n-t+1)(1-X_{(t-1)j})} - 2(\log(n-t+1)+1) \leq 1 + \sqrt{2(n-t+1)(1-h_{n-t+1}(X_{(t-1)j}))},$$

which in turn implies that

$$\begin{aligned} h_{n-t+1}(X_{(t-1)j}) - X_{(t-1)j} &\leq \frac{2\log(n-t+1)+3}{\sqrt{2(n-t+1)}} \left(\sqrt{1-X_{(t-1)j}} + \sqrt{1-h_{n-t+1}(X_{(t-1)j})} \right) \\ &\leq \frac{4\log(n-t+1)+6}{\sqrt{2(n-t+1)}}. \end{aligned}$$

□

Moreover, in this context, the Doob's martingale is simply $M_t = \mathbb{E}[R_n | \mathcal{F}_t]$ for $t = 0, 1, \dots, n$, and hence the associated martingale differences are

$$d_t = M_t - M_{t-1} = \mathbb{E}[R_n | \mathcal{F}_t] - \mathbb{E}[R_n | \mathcal{F}_{t-1}], \quad \text{for all } t \in [n].$$

Assumption 5.1 (Bounded copulas). The joint distribution $F : [0, 1]^p \rightarrow [0, 1]$ satisfies that for any pair $i, j \in [p]$ and $i \neq j$, the bivariate distribution $F_{ij}(w_i, w_j)$ is smooth with density (copula) f_{ij} . Moreover, there exists constant \bar{c} such that

$$f_{ij}(y, z) \leq \bar{c}, \quad \text{for all } i, j \in [p], i \neq j \text{ and all } (y, z) \in [0, 1]^2.$$

In addition, for each $j \in [p]$, there are at most $\theta(n, p)$ coordinates that are correlated with j .

Remark 5.2. The sequential monotone subsequence selection problem is distribution invariant for any continuous distribution. As a result, whenever the marginal distributions are all continuous, it is without loss of generality to assume that all the marginal distributions are the uniform distribution on the unit interval.

When the joint distribution of the uncertainty satisfies Assumption 5.1, we are able to use bootstrap method to simultaneously construct confidence intervals for the total rewards collected by each decision maker, even when p is much larger than n . Such results are summarized in the following theorem.

Theorem 5.4. Suppose that the distribution of uncertainty F satisfied Assumption 5.1. If

$$p \leq \frac{\exp(n^{2a} \log^3 n)}{n^{3/4-3a/2}} \quad \text{and} \quad \theta(n, p) \leq \frac{n^{1/8-3a/4} \log(np)}{\log^2 p \log^{3/2} n}, \quad (5.24)$$

then the simultaneous confidence intervals for the total rewards $R = (R_{n1}, \dots, R_{np})$

$$\prod_{j=1}^p \left[\sqrt{2n} + \sum_{t=1}^n \hat{d}_{tj} - \hat{c}(\alpha)n^{1/4}, \sqrt{2n} + \sum_{t=1}^n \hat{d}_{tj} + \hat{c}(\alpha)n^{1/4} \right]$$

as constructed in Section 5.4 are asymptotic valid. That is,

$$\mathbb{P} \left\{ \sqrt{2n} + \sum_{t=1}^n \hat{d}_{tj} - \hat{c}(\alpha)n^{1/4} \leq R_{nj} \leq \sqrt{2n} + \sum_{t=1}^n \hat{d}_{tj} + \hat{c}(\alpha)n^{1/4} \right\} \rightarrow (1 - \alpha), \quad \text{as } n \rightarrow \infty.$$

The proof of Theorem 5.4 relies on a detailed understanding of each individual MDP in order to verify the conditions required by Theorem 5.2. Such detailed properties of the individual MDP are stated in the following two subsections, after which we prove Theorem 5.4.

5.6.2 Control the Covariance Matrix

The variance of the total reward R_{nj} grows sublinearly, which requires a different normalization than the usual case. Specifically, it is well-known that $\text{Var} \{R_{nj}\} \sim \sqrt{2n}/3$ as $n \rightarrow \infty$ (Arlotto et al., 2015, Theorem 1). Therefore, we scale up the martingale differences by $n^{1/4}$ to make the variance grow linearly in n . That is, we set

$$Z_{tj} \equiv n^{1/4} d_{tj} \quad \text{for all } t \in [n] \text{ and all } j \in [p]. \quad (5.25)$$

As a result, the diagonal entries of the covariance matrix

$$V \equiv \frac{1}{n} \sum_{t=1}^n \mathbb{E}[Z_t Z_t'] = \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbb{E}[\mathbf{d}_t \mathbf{d}_t'] \quad (5.26)$$

would converge to constants as n goes to infinity. We also have the counterpart

$$V_n \equiv \frac{1}{n} \sum_{t=1}^n \mathbb{E}[Z_t Z_t' \mid \mathcal{F}_{t-1}] = \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbb{E}[\mathbf{d}_t \mathbf{d}_t' \mid \mathcal{F}_{t-1}]. \quad (5.27)$$

Under suitable conditions of the uncertainty, we have the convergence in probability $V_n \xrightarrow{p} V$ as $n \rightarrow \infty$ with well controlled convergence rate. This result is summarized in the following proposition. We note that for the tail probability of $\left\| \frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n \right\|_\infty$, the Azuma-Hoeffding inequality would not be powerful enough due to the unusual normalization $n^{-1/2}$ on the right-hand side of (5.27). For instance, if we apply Azuma-Hoeffding inequality to the j -th diagonal entry of $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$ by noticing $|d_{tj}| \leq 1$ for all $t \in [n]$, we would obtain that for any $\lambda > 0$,

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n \right|_{jj} > \lambda \right\} = \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n (d_{tj}^2 - \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]) \right| > \lambda \right\} \leq 2 \exp(-\lambda^2/2),$$

which is too loose for our purpose.

Proposition 5.5 (Convergence in probability of the covariance matrix). *If the uncertainty distribution F satisfies Assumption 5.1, then there exists a constant $c_1 > 0$ depends only on \bar{c} , such that for any fixed $a \in (0, 1/6)$, with probability at least $1 - 3p^2 \exp(-n^{2a})$, we have that*

$$\max_{i,j \in [p]} |V_{n,ij} - V_{ij}| \leq \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}},$$

as well as that

$$\max_{i,j \in [p]} \left| \frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n \right| \leq \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}}.$$

As a result, we also have that

$$\mathbb{E} \left[\max_{i,j \in [p]} |V_n - V|_{ij} \right] \leq \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} + 3p^2 n^{1/2} \exp(-n^{2a}).$$

In order to prove Proposition 5.5, we need some intermediate results. Essentially, we are interested in the convergences in probability

$$\|V_n - V\|_\infty \xrightarrow{p} 0 \quad \text{and} \quad \left\| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right\|_\infty \xrightarrow{p} 0 \quad \text{as } n \rightarrow \infty.$$

To facilitate our analysis, we define another sequence of martingale differences $\Delta_t \in \mathbb{R}^{p \times p}$ for all $t \in [n]$ as

$$\Delta_t \equiv \mathbb{E}[V_n \mid \mathcal{F}_t] - \mathbb{E}[V_n \mid \mathcal{F}_{t-1}], \quad (5.28)$$

so that we have the representation

$$V_n - V = \sum_{t=1}^n \Delta_t. \quad (5.29)$$

The proof of Proposition 5.5 requires two intermediate technical propositions. The first technical proposition is the following.

Proposition 5.6. Suppose the uncertainty distribution F satisfies Assumption 5.1, then for all $j \in [p]$, the j -th diagonal entry of $|V_n - V|$ has the tail probability bound that for any $\lambda > 0$ and $\nu \geq \sqrt{2}(6 + 4 \log n/3)^2/\sqrt{n}$,

$$\mathbb{P}\left\{|V_n - V|_{jj} > \lambda\right\} \leq 2 \exp\left(-\frac{\lambda^2}{2\nu}\right) + \exp\left(-\frac{\left(\frac{n\nu}{(6+4 \log n/3)^2} - \sqrt{2n}\right)^2}{2n}\right).$$

For all $i, j \in [p]$ with $i \neq j$, the (i, j) -entry of $|V_n - V|$ has the tail probability bound that for any $\lambda > 0$ and $\nu \geq 8\sqrt{2}(\bar{c} + 1)^2(2 \log n + 3)^6/(9\sqrt{n})$,

$$\mathbb{P}\left\{|V_n - V|_{ij} > \lambda\right\} \leq 2 \exp\left(-\frac{\lambda^2}{2\nu}\right) + \exp\left(-\frac{\left(\frac{9n\nu}{4(\bar{c}+1)^2(2 \log n+3)^6} - 2\sqrt{2n}\right)^2}{2n}\right).$$

Proof. Step 1. Diagonal entries of $V_n - V$.

From Arlotto et al. (2015, Lemma 24) we know that

$$|(\Delta_t)_{jj}| \leq \frac{6 + \frac{4}{3} \log n}{\sqrt{n}} \mathbb{1}\{G_t^j\}, \quad (5.30)$$

where the event G_t^j is when a selection happens, that is,

$$G_t^j \equiv \{\omega : W_{tj} \in [X_{(t-1)j}, h_{n-t+1}(X_{(t-1)j})]\}.$$

We know from (5.29) that the tail probability of the j -th diagonal entry of $V_n - V$ can be represented as

$$\begin{aligned} \mathbb{P}\left\{|V_n - V|_{jj} > \lambda\right\} &= \mathbb{P}\left\{\left|\sum_{t=1}^n (\Delta_t)_{jj}\right| > \lambda, \sum_{t=1}^n (\Delta_t)_{jj}^2 \leq \nu\right\} + \mathbb{P}\left\{\left|\sum_{t=1}^n (\Delta_t)_{jj}\right| > \lambda, \sum_{t=1}^n (\Delta_t)_{jj}^2 > \nu\right\} \\ &\leq \mathbb{P}\left\{\left|\sum_{t=1}^n (\Delta_t)_{jj}\right| > \lambda, \sum_{t=1}^n (\Delta_t)_{jj}^2 \leq \nu\right\} + \mathbb{P}\left\{\sum_{t=1}^n (\Delta_t)_{jj}^2 > \nu\right\}. \end{aligned} \quad (5.31)$$

To bound the first probability, we apply Belloni and Oliveira (2017, Lemma 9) to obtain that

$$\mathbb{P}\left\{\left|\sum_{t=1}^n (\Delta_t)_{jj}\right| > \lambda, \sum_{t=1}^n (\Delta_t)_{jj}^2 \leq \nu\right\} \leq 2 \exp\left(-\frac{\lambda^2}{2\nu}\right). \quad (5.32)$$

On the other hand, we know from (5.30) that the second probability on the right-hand side of (5.31) is bounded by

$$\mathbb{P}\left\{\sum_{t=1}^n (\Delta_t)_{jj}^2 > \nu\right\} \leq \mathbb{P}\left\{\sum_{t=1}^n \mathbb{1}\{G_t^j\} > \frac{n\nu}{(6 + 4 \log n/3)^2}\right\}.$$

Then we apply Azuma-Hoeffding inequality to obtain that

$$\begin{aligned} &\mathbb{P}\left\{\sum_{t=1}^n \mathbb{1}\{G_t^j\} > \frac{n\nu}{(6 + 4 \log n/3)^2}\right\} \\ &= \mathbb{P}\left\{\sum_{t=1}^n \left\{\mathbb{1}\{G_t^j\} - \mathbb{E}\left[\mathbb{1}\{G_t^j\}\right]\right\} > \frac{n\nu}{(6 + 4 \log n/3)^2} - \sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\{G_t^j\}\right]\right\} \\ &\leq \exp\left(-\frac{\left(\frac{n\nu}{(6 + 4 \log n/3)^2} - \sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\{G_t^j\}\right]\right)^2}{2n}\right). \end{aligned}$$

Another application of Arlotto et al. (2015, Theorem 1, Equation (5)) yields that $\sum_{t=1}^n \mathbb{E}\left[\mathbb{1}\{G_t^j\}\right] \leq \sqrt{2n}$, so if we choose ν such that

$$\frac{n\nu}{(6 + 4 \log n/3)^2} \geq \sqrt{2n},$$

then we obtain an upper bound on the second piece of the right-hand side of (5.31), that is,

$$\mathbb{P} \left\{ \sum_{t=1}^n (\Delta_t)_{jj}^2 > \nu \right\} \leq \exp \left(- \frac{\left(\frac{n\nu}{(6+4\log n/3)^2} - \sqrt{2n} \right)^2}{2n} \right). \quad (5.33)$$

Combining (5.32) and (5.33) together, we finally obtain the tail probability bound on the j -th diagonal entry of $|V_n - V|$ given as

$$\mathbb{P} \left\{ |V_n - V|_{jj} > \lambda \right\} \leq 2 \exp \left(- \frac{\lambda^2}{2\nu} \right) + \exp \left(- \frac{\left(\frac{n\nu}{(6+4\log n/3)^2} - \sqrt{2n} \right)^2}{2n} \right),$$

for all $\lambda > 0$ and all $\nu \geq \sqrt{2/n}(6 + 4\log n/3)^2$.

Step 2. Off-diagonal entries of $V_n - V$.

For $i, j \in [p]$ and $i \neq j$, the (i, j) -element of the martingale difference Δ_t from (5.28) is simply

$$\begin{aligned} (\Delta_t)_{ij} &= \mathbb{E}[(V_n)_{ij} \mid \mathcal{F}_t] - \mathbb{E}[(V_n)_{ij} \mid \mathcal{F}_{t-1}] \\ &= \frac{1}{\sqrt{n}} \sum_{k=t+1}^n \mathbb{E}[d_{ki}d_{kj} \mid \mathcal{F}_t] - \mathbb{E} \left[\frac{1}{\sqrt{n}} \sum_{k=t+1}^n \mathbb{E}[d_{ki}d_{kj} \mid \mathcal{F}_t] \mid \mathcal{F}_{t-1} \right]. \end{aligned} \quad (5.34)$$

If we introduce the shorthand $\omega_{n-t}(X_{ti}, X_{tj}) \equiv \sum_{k=t+1}^n \mathbb{E}[d_{ki}d_{kj} \mid \mathcal{F}_t]$, then this last equation can be re-written as

$$\sqrt{n}(\Delta_t)_{ij} = \omega_{n-t}(X_{ti}, X_{tj}) - \mathbb{E}[\omega_{n-t}(X_{ti}, X_{tj}) \mid \mathcal{F}_{t-1}]. \quad (5.35)$$

Each one of the original martingale differences, d_{ti} , can be represented as the sum of two quantities depending on what happens at period t . Specifically, one has that $d_{ti} = A_{ti} + B_{ti}$ (c.f. Arlotto et al., 2015, Equation (45)), where

$$\begin{aligned} A_{ti} &\equiv (1 + v_{n-t}(W_{ti}) - v_{n-t}(X_{(t-1)i})) \mathbb{1} \{W_{ti} \in [X_{(t-1)i}, h_{n-t+1}(X_{(t-1)i})]\}, \\ B_{ti} &\equiv v_{n-t}(X_{(t-1)i}) - v_{n-t+1}(X_{(t-1)i}). \end{aligned}$$

Since d_{ti} is a martingale difference, we have the identity $\mathbb{E}[A_{ti} \mid \mathcal{F}_{t-1}] = B_{ti}$. As a result,

$$\mathbb{E}[d_{ti}d_{tj} \mid \mathcal{F}_{t-1}] = \mathbb{E}[A_{ti}A_{tj} \mid \mathcal{F}_{t-1}] - B_{ti}B_{tj}.$$

The \mathcal{F}_{t-1} -measurable quantity $B_{ti}B_{tj}$ is easy to bound since for each one of the two pieces, we

have that

$$\begin{aligned}
|B_{ti}| &= |v_{n-t}(X_{(t-1)i}) - v_{n-t+1}(X_{(t-1)i})| \\
&= \left| \int_{X_{(t-1)i}}^{h_{n-t+1}(X_{(t-1)i})} \{1 + v_{n-t}(w) - v_{n-t}(X_{(t-1)i})\} dw \right| \\
&\leq h_{n-t+1}(X_{(t-1)i}) - X_{(t-1)i},
\end{aligned}$$

and from (5.23) we know that this last right-hand side is at most $\frac{4\log(n-t+1)+6}{\sqrt{2(n-t+1)}}$. As a result,

$$|B_{ti}B_{tj}| \leq \frac{2(2\log(n-t+1)+3)^2}{(n-t+1)}. \quad (5.36)$$

On the other hand, we have the integral representation

$$\begin{aligned}
&\mathbb{E}[A_{ti}A_{tj} \mid \mathcal{F}_{t-1}] \\
&= \iint_{\Lambda} \{1 + v_{n-t}(w_1) - v_{n-t}(X_{(t-1)i})\} \{1 + v_{n-t}(w_2) - v_{n-t}(X_{(t-1)j})\} f_{ij}(w_1, w_2) dw_1 dw_2,
\end{aligned}$$

where the double integral is over the region $\Lambda \equiv [X_{(t-1)i}, h_{n-t+1}(X_{(t-1)i})] \times [X_{(t-1)j}, h_{n-t+1}(X_{(t-1)j})]$, and $f_{ij}(\cdot, \cdot)$ is the density of the joint uncertainty faced by the i -th and the j -th decision makers.

Then by Assumption 5.1 we know that

$$\mathbb{E}[A_{ti}A_{tj} \mid \mathcal{F}_{t-1}] \leq \bar{c} (h_{n-t+1}(X_{(t-1)i}) - X_{(t-1)i}) (h_{n-t+1}(X_{(t-1)j}) - X_{(t-1)j}),$$

and if we apply (5.23) once again, we would obtain that

$$\mathbb{E}[A_{ti}A_{tj} \mid \mathcal{F}_{t-1}] \leq \frac{2\bar{c}(2\log(n-t+1)+3)^2}{(n-t+1)}. \quad (5.37)$$

Finally, we assemble the estimations (5.37) and (5.36) together to obtain that

$$\begin{aligned}
|\omega_{n-t}(X_{ti}, X_{tj})| &= \left| \sum_{k=t+1}^n \mathbb{E}[d_{ki}d_{kj} \mid \mathcal{F}_t] \right| = \left| \sum_{k=t+1}^n \mathbb{E}[A_{ki}A_{kj} \mid \mathcal{F}_t] - \sum_{k=t+1}^n \mathbb{E}[B_{ki}B_{kj} \mid \mathcal{F}_t] \right| \\
&= \sum_{k=t+1}^n \frac{2(\bar{c}+1)(2\log(n-k+1)+3)^2}{(n-k+1)} \\
&\leq \frac{2(\bar{c}+1)(2\log(n-t)+3)^3}{3}. \quad (5.38)
\end{aligned}$$

Consider the following three disjoint events that describe what could happen at time period t ,

$$\begin{aligned} G_t^{ij} &\equiv \{\omega : W_{ti} \in [X_{(t-1)i}, h_{n-t+1}(X_{(t-1)i})], W_{tj} \in [X_{(t-1)j}, h_{n-t+1}(X_{(t-1)j})]\} \\ G_t^i &\equiv \{\omega : W_{ti} \in [X_{(t-1)i}, h_{n-t+1}(X_{(t-1)i})], W_{tj} \notin [X_{(t-1)j}, h_{n-t+1}(X_{(t-1)j})]\} \\ G_t^j &\equiv \{\omega : W_{ti} \notin [X_{(t-1)i}, h_{n-t+1}(X_{(t-1)i})], W_{tj} \in [X_{(t-1)j}, h_{n-t+1}(X_{(t-1)j})]\}. \end{aligned}$$

In words, the event G_t^{ij} is that both the i -th decision maker and the j -th decision maker make a selection in time period t ; the event G_t^i is that the i -th decision maker makes a selection in time period t while the j -th decision maker does not; the event G_t^j is that the j -th decision maker makes a selection in time period t while the i -th decision maker does not. Then we can represent $\omega_{n-t}(X_{ti}, X_{tj})$ in terms of the three events as

$$\begin{aligned} \omega_{n-t}(X_{ti}, X_{tj}) &= \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j}) + \{\omega_{n-t}(W_{ti}, W_{tj}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^{ij}\} \\ &\quad + \{\omega_{n-t}(W_{ti}, X_{(t-1)j}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^i\} \\ &\quad + \{\omega_{n-t}(X_{(t-1)i}, W_{tj}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^j\}. \end{aligned}$$

Since the first term on this last right-hand side is \mathcal{F}_{t-1} -measurable, applying this last equation to (5.35) implies that

$$\begin{aligned} \sqrt{n}(\Delta_t)_{ij} &= \{\omega_{n-t}(W_{ti}, W_{tj}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^{ij}\} \\ &\quad - \mathbb{E} \left[\{\omega_{n-t}(W_{ti}, W_{tj}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^{ij}\} \mid \mathcal{F}_{t-1} \right] \\ &\quad + \{\omega_{n-t}(W_{ti}, X_{(t-1)j}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^i\} \\ &\quad - \mathbb{E} \left[\{\omega_{n-t}(W_{ti}, X_{(t-1)j}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^i\} \mid \mathcal{F}_{t-1} \right] \\ &\quad + \{\omega_{n-t}(X_{(t-1)i}, W_{tj}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^j\} \\ &\quad - \mathbb{E} \left[\{\omega_{n-t}(X_{(t-1)i}, W_{tj}) - \omega_{n-t}(X_{(t-1)i}, X_{(t-1)j})\} \mathbb{1}\{G_t^j\} \mid \mathcal{F}_{t-1} \right]. \end{aligned}$$

Next, if we apply the uniform bound (5.38), we would obtain that

$$|(\Delta_t)_{ij}| \leq \frac{2(\bar{c} + 1)(2 \log(n-t) + 3)^3}{3\sqrt{n}} \mathbb{1}\{G_t^{ij} \cup G_t^i \cup G_t^j\}, \quad (5.39)$$

and so the same argument from *Step 1* applies here. Basically, from (5.29) we know that the tail

probability of $|V_n - V|_{ij}$ satisfies that

$$\begin{aligned}
\mathbb{P} \left\{ |V_n - V|_{ij} > \lambda \right\} &= \mathbb{P} \left\{ \left| \sum_{t=1}^n (\Delta_t)_{ij} \right| > \lambda \right\} \\
&= \mathbb{P} \left\{ \left| \sum_{t=1}^n (\Delta_t)_{ij} \right| > \lambda, \sum_{t=1}^n (\Delta_t)_{ij}^2 \leq \nu \right\} + \mathbb{P} \left\{ \left| \sum_{t=1}^n (\Delta_t)_{ij} \right| > \lambda, \sum_{t=1}^n (\Delta_t)_{ij}^2 > \nu \right\} \\
&\leq \mathbb{P} \left\{ \left| \sum_{t=1}^n (\Delta_t)_{ij} \right| > \lambda, \sum_{t=1}^n (\Delta_t)_{ij}^2 \leq \nu \right\} + \mathbb{P} \left\{ \sum_{t=1}^n (\Delta_t)_{ij}^2 > \nu \right\}.
\end{aligned}$$

We apply Belloni and Oliveira (2017, Lemma 9) to the first term of this last right-hand side to obtain that

$$\mathbb{P} \left\{ \left| \sum_{t=1}^n (\Delta_t)_{ij} \right| > \lambda, \sum_{t=1}^n (\Delta_t)_{ij}^2 \leq \nu \right\} \leq 2 \exp \left(-\frac{\lambda^2}{2\nu} \right). \quad (5.40)$$

On the other hand, we apply (5.39) to obtain that

$$\mathbb{P} \left\{ \sum_{t=1}^n (\Delta_t)_{ij}^2 > \nu \right\} \leq \mathbb{P} \left\{ \sum_{t=1}^n \mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} > \frac{9n\nu}{4(\bar{c}+1)^2 (2 \log n + 3)^6} \right\}.$$

Then we apply Azuma-Hoeffding inequality to obtain that

$$\begin{aligned}
&\mathbb{P} \left\{ \sum_{t=1}^n \mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} > \frac{9n\nu}{4(\bar{c}+1)^2 (2 \log n + 3)^6} \right\} \\
&= \mathbb{P} \left\{ \sum_{t=1}^n \mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} - \mathbb{E} \left[\mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} \right] > \frac{9n\nu}{4(\bar{c}+1)^2 (2 \log n + 3)^6} - \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} \right] \right\} \\
&\leq \exp \left(-\frac{\left(\frac{9n\nu}{4(\bar{c}+1)^2 (2 \log n + 3)^6} - \sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} \right] \right)^2}{2n} \right).
\end{aligned}$$

Since one always has

$$\mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} \leq \mathbb{1} \left\{ G_t^i \right\} + \mathbb{1} \left\{ G_t^j \right\},$$

and

$$\sum_{t=1}^n \mathbb{E} \left[\mathbb{1} \left\{ G_t^i \right\} \right] \leq \sqrt{2n},$$

so if we choose ν such that

$$\frac{9n\nu}{4(\bar{c}+1)^2 (2 \log n + 3)^6} \geq 2\sqrt{2n},$$

then we obtain an upper bound

$$\begin{aligned} \mathbb{P} \left\{ \sum_{t=1}^n (\Delta_t)_{ij}^2 > \nu \right\} &\leq \mathbb{P} \left\{ \sum_{t=1}^n \mathbb{1} \left\{ G_t^{ij} \cup G_t^i \cup G_t^j \right\} > \frac{9n\nu}{4(\bar{c}+1)^2 (2\log n + 3)^6} \right\} \\ &\leq \exp \left(- \frac{\left(\frac{9n\nu}{4(\bar{c}+1)^2 (2\log n + 3)^6} - 2\sqrt{2n} \right)^2}{2n} \right). \end{aligned} \quad (5.41)$$

Putting (5.40) and (5.41) together, we obtain that

$$\mathbb{P} \left\{ |V_n - V|_{ij} > \lambda \right\} \leq 2 \exp \left(- \frac{\lambda^2}{2\nu} \right) + \exp \left(- \frac{\left(\frac{9n\nu}{4(\bar{c}+1)^2 (2\log n + 3)^6} - 2\sqrt{2n} \right)^2}{2n} \right),$$

for all $\lambda > 0$ and all $\nu \geq 8\sqrt{2}(\bar{c}+1)^2(2\log n + 3)^6/(9\sqrt{n})$. \square

The second intermediate technical proposition is the following.

Proposition 5.7. For all $\lambda > 0$ and all $\nu \geq \frac{\sqrt{2n}+2(\log n+1)}{3n}$, the j -th diagonal entry of the matrix $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$ has the tail probability bound

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n \right|_{jj} > \lambda \right\} \leq 2 \exp \left[- \frac{\lambda^2}{2\nu} \left(2 - \exp \left(\frac{\lambda}{\sqrt{n\nu}} \right) \right) \right] + \exp \left(- \frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n) \right)^2}{2n \left(6 + \frac{4}{3} \log n \right)} \right).$$

For any pair (i, j) with $i \neq j$, the (i, j) -entry of the matrix $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$ has the tail probability bound

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n \right|_{ij} > \lambda \right\} \leq 2 \exp \left[- \frac{\lambda^2}{2\nu} \left(2 - \exp \left(\frac{\lambda}{\sqrt{n\nu}} \right) \right) \right] + \exp \left(- \frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n) \right)^2}{2n \left(6 + \frac{4}{3} \log n \right)} \right).$$

Proof. Step 1. Diagonal entries of $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$.

From (5.25) and (5.27), we know that the j -th diagonal entry of the matrix $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$ is given by

$$\frac{1}{\sqrt{n}} \sum_{t=1}^n d_{tj}^2 - \mathbb{E} [d_{tj}^2 \mid \mathcal{F}_{t-1}],$$

so with the shorthand $\eta_{tj} \equiv d_{tj}^2 - \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]$ for all $t \in [n]$ and all $j \in [p]$, one can easily verify that $\{\eta_{tj}\}_{t=1}^n$ is a sequence of martingale differences as $\mathbb{E}[\eta_{tj} \mid \mathcal{F}_{t-1}] = 0$. Moreover, for all $j \in [p]$, we also have the tail probability bound

$$\begin{aligned} & \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n d_{tj}^2 - \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] \right| > \lambda \right\} = \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tj} \right| > \lambda \right\} \\ &= \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tj} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] \leq \nu \right\} + \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tj} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \\ &\leq \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tj} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] \leq \nu \right\} + \mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] > \nu \right\}. \end{aligned}$$

As for the first piece of this last right-hand side, note that $|\eta_{tj}| \leq 1$, so we apply Belloni and Oliveira (2017, Lemma 10) to obtain that

$$\begin{aligned} & \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tj} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] \leq \nu \right\} \\ &= \mathbb{P} \left\{ \left| \sum_{t=1}^n \eta_{tj} \right| > \sqrt{n}\lambda, \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] \leq n\nu \right\} \leq 2 \exp \left[-\frac{\lambda^2}{2\nu} \left(2 - \exp\left(\frac{\lambda}{\sqrt{n\nu}}\right) \right) \right]. \end{aligned} \quad (5.42)$$

While for the second piece, from the definition $\eta_{tj} = d_{tj}^2 - \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]$ we know that

$$\mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] = \mathbb{E}[d_{tj}^4 \mid \mathcal{F}_{t-1}] - \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]^2 \leq \mathbb{E}[d_{tj}^4 \mid \mathcal{F}_{t-1}],$$

and since $|d_{tj}| \leq 1$ (see, e.g., Arlotto et al., 2015, Equation (43)), we further have that $\mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] \leq \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}]$. As a result,

$$\begin{aligned} & \mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tj}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \leq \mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \\ &= \mathbb{P} \left\{ \sum_{t=1}^n \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] - \mathbb{E}[d_{tj}^2] > n\nu - \sum_{t=1}^n \mathbb{E}[d_{tj}^2] \right\}. \end{aligned}$$

In the rightmost probability of this last equation, we can once again represent the zero-mean quantity as a sum of martingale differences. Recall the two matrices (5.26) and (5.27), we know that

$$\sum_{t=1}^n \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] - \mathbb{E}[d_{tj}^2] = \sqrt{n}(V_n - V)_{jj} = \sum_{t=1}^n \sqrt{n}(\Delta_t)_{jj}, \quad (5.43)$$

and from (5.30) we further obtain that each martingale difference is bounded by

$$\sqrt{n} |(\Delta_t)_{jj}| \leq \left(6 + \frac{4}{3} \log n\right) \mathbb{1} \{G_t^j\}.$$

Then we apply Azuma-Hoeffding inequality to the right tail of the sum of the martingale differences in (5.43) to obtain that

$$\mathbb{P} \left\{ \sum_{t=1}^n \mathbb{E} [d_{tj}^2 \mid \mathcal{F}_{t-1}] - \mathbb{E} [d_{tj}^2] > n\nu - \sum_{t=1}^n \mathbb{E} [d_{tj}^2] \right\} \leq \exp \left(- \frac{(n\nu - \sum_{t=1}^n \mathbb{E} [d_{tj}^2])^2}{2n \left(6 + \frac{4}{3} \log n\right)} \right).$$

By Arlotto et al. (2015, Theorem 1 and Proposition 20) we know that

$$\sum_{t=1}^n \mathbb{E} [d_{tj}^2] = \text{Var} \left\{ \sum_{t=1}^n d_{tj} \right\} \leq \frac{\sqrt{2n}}{3} + \frac{2}{3}(1 + \log n), \quad (5.44)$$

which in turn implies that if $n\nu \geq \frac{\sqrt{2n}}{3} + \frac{2}{3}(1 + \log n)$, then

$$\mathbb{P} \left\{ \sum_{t=1}^n \mathbb{E} [d_{tj}^2 \mid \mathcal{F}_{t-1}] - \mathbb{E} [d_{tj}^2] > n\nu - \sum_{t=1}^n \mathbb{E} [d_{tj}^2] \right\} \leq \exp \left(- \frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n)\right)^2}{2n \left(6 + \frac{4}{3} \log n\right)} \right), \quad (5.45)$$

and so we finally obtain that

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E} [\eta_{tj}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \leq \exp \left(- \frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n)\right)^2}{2n \left(6 + \frac{4}{3} \log n\right)} \right). \quad (5.46)$$

Putting (5.42) and (5.46) together, we obtain the final bound that for all $t > 0$ and all $\nu \geq \frac{\sqrt{2n} + 2(1 + \log n)}{3n}$,

$$\begin{aligned} & \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n d_{tj}^2 - \mathbb{E} [d_{tj}^2 \mid \mathcal{F}_{t-1}] \right| > \lambda \right\} \\ & \leq 2 \exp \left[- \frac{\lambda^2}{2\nu} \left(2 - \exp \left(\frac{\lambda}{\sqrt{n\nu}} \right) \right) \right] + \exp \left(- \frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n)\right)^2}{2n \left(6 + \frac{4}{3} \log n\right)} \right). \end{aligned}$$

Step 2. Off diagonal entries of $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$.

The (i, j) -element of the matrix $\frac{1}{n} \sum_{t=1}^n Z_t Z_t' - V_n$ is given by

$$\frac{1}{\sqrt{n}} \sum_{t=1}^n d_{ti} d_{tj} - \mathbb{E}[d_{ti} d_{tj} \mid \mathcal{F}_{t-1}],$$

so with the shorthand $\eta_{tij} \equiv d_{ti} d_{tj} - \mathbb{E}[d_{ti} d_{tj} \mid \mathcal{F}_{t-1}]$ for all $t \in [n]$ and all $i, j \in [p]$ with $i \neq j$, one can easily verify that $\{\eta_{tij}\}_{t \in [n]}$ is a sequence of martingale differences as $\mathbb{E}[\eta_{tij} \mid \mathcal{F}_{t-1}] = 0$. Moreover, for all $i, j \in [p]$ with $i \neq j$, we also have the tail probability bound

$$\begin{aligned} & \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n d_{ti} d_{tj} - \mathbb{E}[d_{ti} d_{tj} \mid \mathcal{F}_{t-1}] \right| > \lambda \right\} = \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tij} \right| > \lambda \right\} \\ &= \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tij} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] \leq \nu \right\} + \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tij} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \\ &\leq \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tij} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] \leq \nu \right\} + \mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] > \nu \right\}. \end{aligned}$$

As for the first piece of this last right-hand side, note that $|\eta_{tij}| \leq 1$, so we apply Belloni and Oliveira (2017, Lemma 10) to obtain that

$$\mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n \eta_{tij} \right| > \lambda, \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] \leq \nu \right\} \leq 2 \exp \left[-\frac{\lambda^2}{2\nu} \left(2 - \exp\left(\frac{\lambda}{\sqrt{n\nu}}\right) \right) \right]. \quad (5.47)$$

While for the second piece, from the definition $\eta_{tij} = d_{ti} d_{tj} - \mathbb{E}[d_{ti} d_{tj} \mid \mathcal{F}_{t-1}]$ and the fact that $|d_{ti}| \leq 1$ (see, e.g., Arlotto et al., 2015, Equation (43)), we have the crude bound

$$\mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] = \mathbb{E}[d_{ti}^2 d_{tj}^2 \mid \mathcal{F}_{t-1}] - \mathbb{E}[d_{ti} d_{tj} \mid \mathcal{F}_{t-1}]^2 \leq \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}].$$

As a result,

$$\begin{aligned} & \mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \leq \mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \\ &= \mathbb{P} \left\{ \sum_{t=1}^n \mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] - \mathbb{E}[d_{tj}^2] > n\nu - \sum_{t=1}^n \mathbb{E}[d_{tj}^2] \right\}. \end{aligned}$$

From (5.45), we know that

$$\mathbb{P} \left\{ \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\eta_{tij}^2 \mid \mathcal{F}_{t-1}] > \nu \right\} \leq \exp \left(-\frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n) \right)^2}{2n \left(6 + \frac{4}{3} \log n \right)} \right). \quad (5.48)$$

Lastly, if we put (5.47) and (5.48) together, we obtain that the (i, j) -element of $\frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n$ has the tail bound

$$\begin{aligned} \mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right|_{ij} > \lambda \right\} &= \mathbb{P} \left\{ \left| \frac{1}{\sqrt{n}} \sum_{t=1}^n d_{ti} d_{tj} - \mathbb{E}[d_{ti} d_{tj} \mid \mathcal{F}_{t-1}] \right| > \lambda \right\} \\ &\leq 2 \exp \left[-\frac{\lambda^2}{2\nu} \left(2 - \exp\left(\frac{\lambda}{\sqrt{n\nu}}\right) \right) \right] + \exp \left(-\frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n) \right)^2}{2n \left(6 + \frac{4}{3} \log n \right)} \right), \end{aligned}$$

so the proof is complete. \square

Now we are ready to prove Proposition 5.5.

Proof. (Proof of Proposition 5.5.) We start with the convergence $\|V_n - V\|_\infty \xrightarrow{p} 0$. From Proposition 5.6 we know that for the j -th diagonal entry of $V_n - V$, its tail probability satisfies that

$$\mathbb{P} \left\{ |V_n - V|_{jj} > \lambda \right\} \leq 2 \exp \left(-\frac{\lambda^2}{2\nu} \right) + \exp \left(-\frac{\left(\frac{n\nu}{(6+4\log n/3)^2} - \sqrt{2n} \right)^2}{2n} \right),$$

so if we fix some $a \in (0, 1/6)$ and choose

$$\lambda = c_1 n^{3a/2-1/4} \log^3 n \quad \text{and} \quad \nu = \frac{c_1 n^{a-1/2} \log^6 n}{2}, \quad (5.49)$$

then we can choose proper constant c_1 to obtain that

$$\mathbb{P} \left\{ |V_n - V|_{jj} > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \leq 3 \exp(-n^{2a}).$$

Similarly, for the (i, j) -entry of $V_n - V$ with $i \neq j$, we know from Proposition 5.6 that its tail probability satisfies that

$$\mathbb{P} \left\{ |V_n - V|_{ij} > \lambda \right\} \leq 2 \exp \left(-\frac{\lambda^2}{2\nu} \right) + \exp \left(-\frac{\left(\frac{9n\nu}{4(\bar{c}+1)^2(2\log n+3)^6} - 2\sqrt{2n} \right)^2}{2n} \right).$$

With the same choice of λ and ν , we can choose proper constant c_1 that depends only on \bar{c} to obtain that

$$\mathbb{P} \left\{ |V_n - V|_{ij} > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \leq 3 \exp(-n^{2a}).$$

Therefore, we have the union bound that

$$\begin{aligned} \mathbb{P} \left\{ \|V_n - V\|_\infty > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} &\leq \sum_{j=1}^n \mathbb{P} \left\{ |V_n - V|_{jj} > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} + \sum_{i \neq j} \mathbb{P} \left\{ |V_n - V|_{ij} > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \\ &\leq 3p^2 \exp(-n^{2a}). \end{aligned}$$

Next, we prove that $\left\| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right\|_\infty \xrightarrow{p} 0$. From Proposition 5.7 we know that for the (i, j) -entry of $\frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n$, its tail probability satisfies that

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right|_{ij} > \lambda \right\} \leq 2 \exp \left[-\frac{\lambda^2}{2\nu} \left(2 - \exp\left(\frac{\lambda}{\sqrt{n\nu}}\right) \right) \right] + \exp \left(-\frac{\left(n\nu - \frac{\sqrt{2n}}{3} - \frac{2}{3}(1 + \log n) \right)^2}{2n \left(6 + \frac{4}{3} \log n \right)} \right).$$

If we choose the same λ and ν as in (5.49), then with proper choice of c_1 we obtain that

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right|_{ij} > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \leq 3 \exp(-n^{2a}).$$

As a result, we have the union bound

$$\mathbb{P} \left\{ \left\| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right\|_\infty > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \leq \sum_{i,j \in [p]} \mathbb{P} \left\{ \left| \frac{1}{n} \sum_{t=1}^n Z_t Z'_t - V_n \right|_{ij} > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \leq 3p^2 \exp(-n^{2a}).$$

Lastly, we prove the upper bound on $\mathbb{E}[\|V_n - V\|_\infty]$. Since we have already proved that

$$\mathbb{P} \left\{ \|V_n - V\|_\infty > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \leq 3p^2 \exp(-n^{2a}),$$

we can bound $\mathbb{E}[\|V_n - V\|_\infty]$ from above by

$$\begin{aligned} \mathbb{E}[\|V_n - V\|_\infty] &\leq \mathbb{P} \left\{ \|V_n - V\|_\infty \leq \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} \frac{\log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} + \mathbb{P} \left\{ \|V_n - V\|_\infty > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} (\|V_n\|_\infty + \|V\|_\infty) \\ &\leq \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} + \mathbb{P} \left\{ \|V_n - V\|_\infty > \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} \right\} (\|V_n\|_\infty + \|V\|_\infty). \end{aligned}$$

Recall the matrices V and V_n from (5.26) and (5.27) receptively, and recall that $|d_{tj}| \leq 1$, we have the crude upper bound $\|V\|_\infty \leq \sqrt{n}$ and $\|V_n\|_\infty \leq \sqrt{n}$. Therefore, with some algebra, we finally obtain that

$$\mathbb{E}[\|V_n - V\|_\infty] \leq \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}} + 3p^2 n^{1/2} \exp(-n^{2a}),$$

and the proof is complete. \square

5.6.3 Verification of Other Conditions

We start with the boundedness condition on $\|V\|_\infty$ in Condition (ii) of Theorem 5.2.

Proposition 5.8 (Boundedness of covariance matrix entries). *The covariance matrix V defined in (5.26) has the uniform bound on all of its entries given as*

$$\max_{i,j \in [p]} |V_{ij}| \leq \frac{\sqrt{2}}{3} + \frac{2(1 + \log n)}{3\sqrt{n}}.$$

Proof. We start with the diagonal entries of V . From the definition of V in (5.26) and the upper bound (5.44), we know that for all $j \in [p]$,

$$V_{jj} = \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbb{E}[d_{tj}^2] = \frac{1}{\sqrt{n}} \mathbb{E} \left[\left(\sum_{t=1}^n d_{tj} \right)^2 \right] = \frac{1}{\sqrt{n}} \text{Var} \{R_{nj}\} \leq \frac{\sqrt{2}}{3} + \frac{2(1 + \log n)}{3\sqrt{n}}.$$

For the off-diagonal entries of V , we know that for each pair $i \neq j$,

$$V_{ij} = \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbb{E}[d_{ti}d_{tj}] = \frac{1}{\sqrt{n}} \mathbb{E} \left[\left(\sum_{t=1}^n d_{ti} \right) \left(\sum_{t=1}^n d_{tj} \right) \right] = \frac{1}{\sqrt{n}} \text{Cov} \{R_{ni}, R_{nj}\} \leq \frac{\sqrt{2}}{3} + \frac{2(1 + \log n)}{3\sqrt{n}},$$

completing the proof. \square

Next, we verify the other conditions in Theorem 5.2. Condition (i) of Theorem 5.2 is verified in the following proposition.

Proposition 5.9 (Boundedness of third moments). *The martingale differences $\{\mathbf{d}_t : t \in [n]\}$ satisfy the moment bound*

$$\sum_{t=1}^n \mathbb{E} \left[\|\mathbf{d}_t\|_\infty^3 \right] \leq 4\sqrt{2n} \log n.$$

Proof. Since we always have that $\|\mathbf{d}_t\|_\infty \leq 1$ (see, e.g., Arlotto et al., 2015, (43)), it suffices for us to work with $\|\mathbf{d}_t\|_\infty^2$. From Arlotto et al. (2015, Equation (41), (46)), we know that

$$\mathbb{E}[d_{tj}^2 \mid \mathcal{F}_{t-1}] \leq \int_{X_{(t-1)j}}^{h_{n-t+1}(X_{(t-1)j})} \{1 + v_{n-t}(w) - v_{n-t}(X_{(t-1)j})\}^2 dw \leq h_{n-t+1}(X_{(t-1)j}) - X_{(t-1)j}.$$

As a result,

$$\mathbb{E}[\|\mathbf{d}_t\|_\infty^2] = \mathbb{E}[\mathbb{E}[\|\mathbf{d}_t\|_\infty^2 \mid \mathcal{F}_{t-1}]] \leq \mathbb{E} \left[\max_{j \in [p]} (h_{n-t+1}(X_{(t-1)j}) - X_{(t-1)j}) \right].$$

Then we apply (5.23) to obtain that

$$\mathbb{E} \left[\|\mathbf{d}_t\|_\infty^2 \right] \leq \frac{4 \log(n-t+1) + 6}{\sqrt{2(n-t+1)}}.$$

Summing up this last inequality yields the higher moment bound

$$\sum_{t=1}^n \mathbb{E} \left[\|\mathbf{d}_t\|_\infty^3 \right] \leq \sum_{t=1}^n \mathbb{E} \left[\|\mathbf{d}_t\|_\infty^2 \right] \leq \sum_{t=1}^n \frac{4 \log(n-t+1) + 6}{\sqrt{2(n-t+1)}} \leq 4\sqrt{2n} \log n,$$

so Condition (i) of Theorem 5.2 is satisfied. \square

The positive definiteness in Condition (ii) of Theorem 5.2 is verified in the following proposition.

Proposition 5.10 (Positive semidefinite condition). When the uncertainty distribution F satisfies Assumption 5.1, there exists constant c_1 that depends only on \bar{c} such that, if δ_n satisfies that

$$\delta_n \geq \frac{c_1 \theta(n, p) \log^2 p \log^3 n}{n^{1/4-3a/2}},$$

then we have that with probability at least $1 - 3p^2 \exp(-n^{2a})$,

$$V_n \leq (1 + \delta_n / \log^2 p) V.$$

Proof. We verify that if the distribution F of the uncertainty satisfies Assumption 5.1, then with high probability, the matrix $\{(1 + \delta_n / \log^2 p) V - V_n\}$ is diagonally dominant, and hence we obtain that

$$V_n \leq (1 + \delta_n / \log^2 p) V.$$

From Proposition 5.5 we know that with probability at least $1 - 3p^2 \exp(-n^{2a})$ for some fixed $a \in (0, 1/6)$,

$$\|V_n - V\|_\infty \leq \frac{c_1 \log^3 n}{n^{\frac{1}{4}-\frac{3a}{2}}}.$$

Since $V_{jj} \geq \sqrt{2}/3 - 2/\sqrt{n}$ (cf. Arlotto et al., 2015, Proposition 19), the diagonal entries of the matrix $\{(1 + \delta_n / \log^2 p) V - V_n\}$ satisfy that with probability at least $1 - 3p^2 \exp(-n^{2a})$,

$$\{(1 + \delta_n / \log^2 p) V - V_n\}_{jj} \geq \frac{\delta_n}{\log^2 p} V_{jj} - \|V_n - V\|_\infty \geq \frac{\sqrt{2}\delta_n}{3\log^2 p} - \frac{c_1 \log^3 n}{n^{\frac{1}{4}-\frac{3a}{2}}}.$$

On the other hand, the off-diagonal entries of the matrix $\{(1 + \delta_n/\log^2 p) V - V_n\}$ has the upper bound that for $i, j \in [p]$ with $i \neq j$,

$$|(1 + \delta_n/\log^2 p) V - V_n|_{ij} \leq \frac{\delta_n}{\log^2 p} |V_{ij}| + \|V_n - V\|_\infty \leq O\left(\frac{\delta_n \log^3 n}{\log^2 p \sqrt{n}}\right) + \frac{c_1 \log^3 n}{n^{\frac{1}{4} - \frac{3a}{2}}}.$$

Combining the last two inequalities together, we obtain that for all $j \in [p]$,

$$\frac{\{(1 + \delta_n/\log^2 p) V - V_n\}_{jj}}{|(1 + \delta_n/\log^2 p) V - V_n|_{ij}} \geq \frac{\delta_n n^{\frac{1}{4} - \frac{3a}{2}}}{c_1 \log^2 p \log^3 n} \quad \text{for all } i \in [p], i \neq j.$$

Hence if the uncertainty distribution F satisfies Assumption 5.1, then for each $j \in [p]$, there are at most $\theta(n, p)$ non-zero entries in each row of the matrix $\{(1 + \delta_n/\log^2 p) V - V_n\}$. The choice of δ_n guarantees that

$$\frac{\delta_n n^{\frac{1}{4} - \frac{3a}{2}}}{c_1 \log^2 p \log^3 n} \geq \theta(n, p),$$

so the matrix $\{(1 + \delta_n/\log^2 p) V - V_n\}$ is diagonally dominant, and hence

$$V_n \leq (1 + \delta_n/\log^2 p) V.$$

□

Proof of Theorem 5.4

With all the conditions verified, we are ready to prove Theorem 5.4.

Proof. (Proof of Theorem 5.4.) We verify the conditions required in Theorem 5.2. We set the parameters

$$\psi_n = 3p^2 \exp(-n^{2a}) \quad \text{and} \quad \delta_n = \frac{\log^{3/2} n \log(np)}{n^{1/8 - 3a/4}}.$$

Then with the parameters satisfying (5.24), one can easily check that: Condition (i) is verified by Proposition 5.9; Condition (ii) is verified by Proposition 5.5, Proposition 5.8, and Proposition 5.10.

The first half of Condition (iii) is verified by Proposition 5.5. Since we can directly compute the martingale differences \mathbf{d}_t with observations, in the second half of Condition (iii) in Theorem 5.2, the estimator \hat{Z}_{kj} is the same as the actual Z_{kj} . As a result, the second half of Condition (iii) is automatically satisfied. □

5.7 Concluding Remarks

In this chapter, we study a general framework of finite-horizon Markov decision problems, and focus on the specific statistical inference problems associated with the optimal and near-optimal policies of such Markov decision problems. The two inferences we focus on are simultaneous hypothesis testing and constructing simultaneous confidence intervals for the total rewards collected by each decision maker. Such inference problems are studied in the high-dimensional regime in which the traditional central limit theorems fail. In particular, we study two different types of MDPs: one with regenerative states and the other with absorbing states. Suitable conditions are given to guarantee asymptotic validity of bootstrap based simultaneous hypothesis testing and simultaneous confidence interval constructions for such two types of MDPs. The key insight from the two applications is that, when one has a comprehensive understanding about the Markov decision problems at hand, it is possible to monitor the implementations of policies of many decision makers simultaneously without knowledge of the correlation structure of the uncertainty they face.

Chapter 6

Conclusion

In this dissertation, we have studied several sequential decision problems under uncertainty. The focus is the total reward collected by different sequential policies. We have analyzed such object from different perspectives. Besides the standard perspective—the expected total reward, we have also analyzed the variance and the limiting distribution of the total reward collected by different sequential policies for certain problems. In addition, we use bootstrap procedure to construct simultaneous confidence intervals for the total rewards collected by each decision maker of many sequential decision problems that are run in parallel. The analyses have both theoretical and practical insights for decision makers facing related problems. Specifically, from a theoretical perspective, our analyses calls for careful comparison when it comes to distinguishing different policies among the class of asymptotically optimal policies for certain problems. The analyses also indicate different risks involved among all asymptotically optimal policies. From a practical perspective, our analyses can be applied to the setting of many sequential decision problems that are correlated and are faced in parallel by self-interested decision makers. The proposed bootstrap procedure can be used for simultaneously testing whether each decision maker is implementing designated policy through the observations of the rewards they collect over time. The tools and techniques developed in this dissertation could potentially find different applications in other related problems.

Bibliography

- Andrews, Donald W. K. 2004. The block–block bootstrap: improved asymptotic refinements. *Econometrica* 72(3) 673–700.
- Arlotto, Alessandro, Itai Gurvich. 2019. Uniformly bounded regret in the multi-secretary problem. *Stochastic Systems* 9(3) 231–260.
- Arlotto, Alessandro, Elchanan Mossel, J. Michael Steele. 2016. Quickest online selection of an increasing subsequence of specified size. *Random Struct. Algor.* 49(2) 235–252.
- Arlotto, Alessandro, Vinh V. Nguyen, J. Michael Steele. 2015. Optimal online selection of a monotone subsequence: a central limit theorem. *Stochastic Process. Appl.* 125(9) 3596–3622.
- Arlotto, Alessandro, J. Michael Steele. 2011. Optimal sequential selection of a unimodal subsequence of a random sequence. *Combin. Probab. Comput.* 20(6) 799–814.
- Arlotto, Alessandro, J. Michael Steele. 2016. A central limit theorem for temporally nonhomogenous Markov chains with applications to dynamic programming. *Math. Oper. Res.* 41(4) 1448–1468.
- Arlotto, Alessandro, Yehua Wei, Xinchang Xie. 2018. An adaptive $O(\log n)$ -optimal policy for the online selection of a monotone subsequence from a random sample. *Random Struct. Algor.* 52(1) 41–53.
- Arlotto, Alessandro, Xinchang Xie. 2020a. Logarithmic regret in the dynamic and stochastic knapsack problem with equal rewards. *Stochastic Systems* forthcoming.
- Arlotto, Alessandro, Xinchang Xie. 2020b. Sequential policies and the distributions of their total rewards in dynamic and stochastic knapsack problems. *Working paper – Duke University* .
- Balseiro, Santiago R., David B. Brown. 2019. Approximations to stochastic dynamic programs via information relaxation duality. *Oper. Res.* 67(2) 577–597.
- Baruah, Sanjoy, Jayant Haritsa, Nitin Sharma. 1994. On-line scheduling to maximize task completions. *1994 Proceedings Real-Time Systems Symposium*. IEEE, San Juan, PR, 228–236.
- Baryshnikov, Yuliy M., Alexander V. Gnedin. 2000. Sequential selection of an increasing sequence from a multidimensional random sample. *Ann. Appl. Probab.* 10(1) 258–267.
- Bellman, Richard. 1957. *Dynamic programming*. Princeton University Press, Princeton, N. J.
- Belloni, Alexandre, Victor Chernozhukov, Denis Chetverikov, Christian Bailey Hansen, Kengo Kato. 2018. High-dimensional econometrics and regularized GMM. *ArXiv e-print 1806.01888* .
- Belloni, Alexandre, Roberto I. Oliveira. 2017. Approximate group context tree. *Ann. Statist.* 45(1) 355–385.
- Belloni, Alexandre, Roberto I. Oliveira. 2018. A high dimensional central limit theorem for martingales, with applications to context tree models. *ArXiv e-print 1809.02741* .
- Belloni, Alexandre, Xinchang Xie. 2020. Data-driven monitoring the optimality of policies across many markov decision problems. *Working paper – Duke University* .

- Bernstein, Serge. 1927. Sur l'extension du théorème limite du calcul des probabilités aux sommes de quantités dépendantes. *Math. Ann.* 97(1) 1–59.
- Bertsekas, Dimitri P. 2005. *Dynamic Programming and Optimal Control*, vol. I,II. Athena Scientific, Belmont, MA.
- Bhalgat, Anand, Ashish Goel, Sanjeev Khanna. 2011. Improved approximation results for stochastic knapsack problems. *Proceedings of the Twenty-Second Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, Philadelphia, PA, 1647–1665.
- Blado, Daniel, Weihong Hu, Alejandro Toriello. 2016. Semi-infinite relaxations for the dynamic knapsack problem with stochastic item sizes. *SIAM J. Optim.* 26(3) 1625–1648.
- Blado, Daniel, Alejandro Toriello. 2019. Relaxation analysis for the dynamic knapsack problem with stochastic item sizes. *SIAM J. Optim.* 29(1) 1–30.
- Boshuizen, Frans A, Robert P Kertz. 1999. Smallest-fit selection of random sizes under a sum constraint: weak convergence and moment comparisons. *Adv. in Appl. Probab.* 31(1) 178–198.
- Bradley, Richard C. 2005. Basic properties of strong mixing conditions. A survey and some open questions. *Probability Surveys* 2 107–144.
- Bruss, F. Thomas, Freddy Delbaen. 2001. Optimal rules for the sequential selection of monotone subsequences of maximum expected length. *Stochastic Process. Appl.* 96(2) 313–342.
- Bruss, F. Thomas, Freddy Delbaen. 2004. A central limit theorem for the optimal selection process for monotone subsequences of maximum expected length. *Stochastic Process. Appl.* 114(2) 287–311.
- Bruss, F. Thomas, Mitia Duerinckx. 2015. Resource dependent branching processes and the envelope of societies. *Ann. Appl. Probab.* 25(1) 324–372.
- Bruss, F. Thomas, James B. Robertson. 1991. “Wald’s lemma” for sums of order statistics of i.i.d. random variables. *Adv. in Appl. Probab.* 23(3) 612–623.
- Bumpensanti, Pornpawee, He Wang. 2020. A re-solving heuristic with uniformly bounded loss for network revenue management. *Management Sci.* forthcoming.
- Carraway, Robert L., Robert L. Schmidt, Lawrence R. Weatherford. 1993. An algorithm for maximizing target achievement in the stochastic knapsack problem with normal returns. *Naval Res. Logist.* 40(2) 161–173.
- Cayley, Arthur. 1875. Mathematical questions and their solutions. *Educational Times* 22 18–19, see *The Collected Mathematical Papers of Arthur Cayley*, 10, 587–588, (1986). Cambridge University Press, Cambridge.
- Chen, Lijian, Tito Homem-de Mello. 2010. Re-solving stochastic programming models for airline revenue management. *Ann. Oper. Res.* 177(1) 91–114.
- Chen, Xiaohong, Qi-Man Shao, Wei Biao Wu, Lihu Xu. 2016. Self-normalized Cramér-type moderate deviations under dependence. *Ann. Statist.* 44(4) 1593–1617.
- Chernozhukov, Victor, Denis Chetverikov, Kengo Kato. 2013. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann. Statist.* 41(6) 2786–2819.

- Chernozhukov, Victor, Denis Chetverikov, Kengo Kato. 2015. Comparison and anti-concentration bounds for maxima of Gaussian random vectors. *Probability Theory and Related Fields* 162(1-2) 47–70.
- Chernozhukov, Victor, Denis Chetverikov, Kengo Kato. 2017. Central limit theorems and bootstrap in high dimensions. *Ann. Probab.* 45(4) 2309–2352.
- Chernozhukov, Victor, Denis Chetverikov, Kengo Kato. 2019. Inference on causal and structural parameters using many moment inequalities. *Rev. Econ. Stud.* 86(5) 1867–1900.
- Coffman, E. G., Jr., L. Flatto, R. R. Weber. 1987. Optimal selection of stochastic intervals under a sum constraint. *Adv. in Appl. Probab.* 19(2) 454–473.
- Cooper, William L. 2002. Asymptotic behavior of an allocation policy for revenue management. *Oper. Res.* 50(4) 720–727.
- Dantzig, George B. 1957. Discrete-variable extremum problems. *Oper. Res.* 5(2) 266–277.
- Dean, Brian C., Michel X. Goemans, Jan Vondrák. 2004. Approximating the stochastic knapsack problem: the benefit of adaptivity. *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science*. IEEE Press, Piscataway, NJ, 208–217. doi:10.1109/FOCS.2004.15.
- Dean, Brian C., Michel X. Goemans, Jan Vondrák. 2005. Adaptivity and approximation for stochastic packing problems. *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*. ACM, New York, NY, 395–404.
- Dean, Brian C., Michel X. Goemans, Jan Vondrák. 2008. Approximating the stochastic knapsack problem: the benefit of adaptivity. *Math. Oper. Res.* 33(4) 945–964.
- Derman, C., G. J. Lieberman, S. M. Ross. 1975. A stochastic sequential allocation model. *Oper. Res.* 23(6) 1120–1130.
- Derman, C., G. J. Lieberman, S. M. Ross. 1978. A renewal decision problem. *Management Sci.* 24(5) 554–561.
- Dobrushin, Roland L. 1956a. Central limit theorem for nonstationary Markov chains. I. *Theory of Probability & Its Applications* 1(1) 65–80.
- Dobrushin, Roland L'vovich. 1956b. Central limit theorem for nonstationary Markov chains. II. *Theory of Probability & Its Applications* 1(4) 329–383.
- Fan, Jianqing, Qiwei Yao. 2003. *Nonlinear Time Series: Nonparametric and Parametric Methods*. Springer.
- Fryzlewicz, Piotr, Suhasini Subba Rao. 2011. Mixing properties of ARCH and time-varying ARCH processes. *Bernoulli* 17(1) 320–346.
- Gallego, Guillermo, Garrett van Ryzin. 1994. Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Management Sci.* 40(8) 999–1020.
- Gallego, Guillermo, Garrett van Ryzin. 1997. A multiproduct dynamic pricing problem and its applications to network yield management. *Oper. Res.* 45(1) 24–41.
- Gnedin, Alexander, Amirlan Seksenbayev. 2019. Asymptotics and renewal approximation in the online selection of increasing subsequence. *ArXiv e-print 1904.11213*.

- Gnedin, Alexander V. 1999. Sequential selection of an increasing subsequence from a sample of random size. *J. Appl. Probab.* 36(4) 1074–1085.
- Gupta, Anupam, Ravishankar Krishnaswamy, Marco Molinaro, R. Ravi. 2011. Approximation algorithms for correlated knapsacks and non-martingale bandits. *2011 IEEE 52nd Annual Symposium on Foundations of Computer Science—FOCS 2011*. IEEE Computer Soc., Los Alamitos, CA, 827–836.
- Henig, Mordechai I. 1990. Risk criteria in a stochastic knapsack problem. *Oper. Res.* 38(5) 820–825.
- Iglehart, Donald L. 1963. Optimality of (s, S) policies in the infinite horizon dynamic inventory problem. *Management Sci.* 9(2) 259–267.
- Ilhan, Taylan, Seyed M. R. Iravani, Mark S. Daskin. 2011. TECHNICAL NOTE—The adaptive knapsack problem with stochastic rewards. *Oper. Res.* 59(1) 242–248.
- Jarner, Søren F., Gareth O. Roberts. 2002. Polynomial convergence rates of Markov chains. *Ann. Appl. Probab.* 12(1) 224–247.
- Jasin, Stefanus, Sunil Kumar. 2012. A re-solving heuristic with bounded revenue loss for network revenue management with customer choice. *Math. Oper. Res.* 37(2) 313–345.
- Jasin, Stefanus, Sunil Kumar. 2013. Analysis of deterministic LP-based booking limit and bid price controls for revenue management. *Oper. Res.* 61(6) 1312–1320.
- Jones, Galin L. 2004. On the Markov chain central limit theorem. *Probability Surveys* 1 299–320.
- Kellerer, Hans, Ulrich Pferschy, David Pisinger. 2004. *Knapsack problems*. Springer-Verlag, Berlin.
- Kleinberg, Robert. 2005. A multiple-choice secretary algorithm with applications to online auctions. *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms*. ACM, New York, NY, 630–631.
- Kleywegt, Anton J., Jason D. Papastavrou. 1998. The dynamic and stochastic knapsack problem. *Oper. Res.* 46(1) 17–35.
- Kleywegt, Anton J., Jason D. Papastavrou. 2001. The dynamic and stochastic knapsack problem with random sized items. *Oper. Res.* 49(1) 26–41.
- Lahiri, Soumendra Nath. 2003. *Resampling methods for dependent data*. Springer Series in Statistics, Springer-Verlag, New York.
- Li, Jian, Wen Yuan. 2013. Stochastic combinatorial optimization via Poisson approximation. *STOC'13—Proceedings of the 2013 ACM Symposium on Theory of Computing*. ACM, New York, NY, 971–980.
- Logan, B. F., L. A. Shepp. 1977. A variational problem for random Young tableaux. *Advances in Math.* 26(2) 206–222.
- Lueker, George S. 1998. Average-case analysis of off-line and on-line knapsack problems. *J. Algorithms* 29(2) 277–305.
- Ma, Will. 2018. Improvements and generalizations of stochastic knapsack and Markovian bandits approximation algorithms. *Math. Oper. Res.* 43(3) 789–812.

- Mandl, Petr, Monika Laušmanová. 1991. Two extensions of asymptotic methods in controlled Markov chains. *Annals of Operations Research* 28(1) 67–79.
- Marchetti-Spaccamela, A., C. Vercellis. 1995. Stochastic on-line knapsack problems. *Math. Program.* 68(1, Ser. A) 73–104.
- Martello, Silvano, Paolo Toth. 1990. *Knapsack problems*. Wiley-Interscience Series in Discrete Mathematics and Optimization, John Wiley & Sons Ltd., Chichester.
- Maxwell, Michael, Michael Woodroffe. 2000. Central limit theorems for additive functionals of Markov chains. *Ann. Probab.* 28(2) 713–724.
- McLeish, D. L. 1974. Dependent central limit theorems and invariance principles. *Ann. Probab.* 2(4) 620–628.
- Mendoza-Pérez, Armando F. 2008. Asymptotic normality of average cost Markov control processes. *Morfismos* 12(2) 33–52.
- Mendoza-Pérez, Armando F., Onésimo Hernández-Lerma. 2010. Asymptotic normality of discrete-time Markov control processes. *J. Appl. Probab.* 47(3) 778–795.
- Merzifonluoğlu, Yasemin, Joseph Geunes, H. Edwin Romeijn. 2012. The static stochastic knapsack problem with normally distributed item sizes. *Math. Program.* 134(2, Ser. A) 459–489.
- Meyn, Sean, Richard L. Tweedie. 2009. *Markov Chains and Stochastic Stability*. 2nd ed. Cambridge University Press, Cambridge.
- Moser, Leo. 1956. On a problem of Cayley. *Scripta Mathematica* 22 289–292.
- Nagaev, S.V. 1976. An estimate of the remainder term in the multidimensional central limit theorem. *Proceedings of the Third Japan—USSR Symposium on Probability Theory*. Springer, 419–438.
- Nakai, Toru. 1986. An optimal selection problem for a sequence with a random number of applicants per period. *Oper. Res.* 34(3) 478–485.
- Paparoditis, Efstathios, Dimitris N. Politis. 2001. Tapered block bootstrap. *Biometrika* 88(4) 1105–1119.
- Paparoditis, Efstathios, Dimitris N. Politis. 2002. The tapered block bootstrap for general statistics from stationary sequences. *Econom. J.* 5(1) 131–148.
- Papastavrou, Jason D., Srikanth Rajagopalan, Anton J. Kleywegt. 1996. The dynamic and stochastic knapsack problem with deadlines. *Management Sci.* 42(12) 1706–1718.
- Peligrad, Magda. 1985. An invariance principle for φ -mixing sequences. *Ann. Probab.* 13(4) 1304–1313.
- Peligrad, Magda. 1990. On Ibragimov–Iosifescu conjecture for φ -mixing sequences. *Stochastic Process. Appl.* 35(2) 293–308.
- Peligrad, Magda. 2012. Central limit theorem for triangular arrays of non-homogeneous Markov chains. *Probability Theory and Related Fields* 154(3–4) 409–428.
- Peng, Peichao, J. Michael Steele. 2016. Sequential selection of a monotone subsequence from a random permutation. *Proc. Amer. Math. Soc.* 144(11) 4973–4982.

- Powell, Warren B. 2011. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. 2nd ed. Wiley Series in Probability and Statistics, John Wiley & Sons, Inc., Hoboken, NJ.
- Prastacos, Gregory P. 1983. Optimal sequential investment decisions under conditions of uncertainty. *Management Sci.* 29(1) 118–134.
- Puterman, Martin L. 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics, John Wiley & Sons, Inc., New York, a Wiley-Interscience Publication.
- Reiman, Martin I., Qiong Wang. 2008. An asymptotically optimal policy for a quantity-based network revenue management problem. *Math. Oper. Res.* 33(2) 257–282.
- Rhee, WanSoo, Michel Talagrand. 1991. A note on the selection of random variables under a sum constraint. *J. Appl. Probab.* 28(4) 919–923.
- Samuels, Stephen M., J. Michael Steele. 1981. Optimal sequential selection of a monotone sequence from a random sample. *Ann. Probab.* 9(6) 937–947.
- Scarf, Herbert. 1959. The optimality of (S, s) policies in the dynamic inventory problem. Kenneth Joseph Arrow, Patrick Suppes, Samuel Karlin, eds., *Mathematical Methods in the Social Sciences*. Stanford University Press, Stanford, California, 196–202.
- Seksenbayev, Amirlan. 2018. Refined asymptotics in the online selection of an increasing subsequence. *ArXiv e-print 1808.06300*.
- Sklar, M. 1959. Fonctions de répartition à n dimensions et leurs marges. *Publ. Inst. Statist. Univ. Paris* 8 229–231.
- Steele, J. Michael. 2016. The Bruss-Robertson inequality: elaborations, extensions, and applications. *Math. Appl. (Warsaw)* 44(1) 3–16.
- Talluri, Kalyan, Garrett van Ryzin. 1998. An analysis of bid-price controls for network revenue management. *Management Sci.* 44(11-part-1) 1577–1593.
- Talluri, Kalyan T., Garrett J. van Ryzin. 2004. *The theory and practice of revenue management*. International Series in Operations Research & Management Science, 68, Kluwer Academic Publishers, Boston, MA.
- Utev, Sergei Aleksandrovich. 1991. On the central limit theorem for φ -mixing arrays of random variables. *Theory of Probability & Its Applications* 35(1) 131–139.
- Vera, Alberto, Siddhartha Banerjee. 2019. The Bayesian prophet: A low-regret framework for online decision making. *SIGMETRICS Perform. Eval. Rev.* 47(1) 81–82.
- Veršik, A. M., S. V. Kerov. 1977. Asymptotic behavior of the Plancherel measure of the symmetric group and the limit form of Young tableaux. *Dokl. Akad. Nauk SSSR* 233(6) 1024–1027.
- White, Douglas J. 1993. A survey of applications of Markov decision processes. *Journal of the Operational Research Society* 44(11) 1073–1096.
- Wu, Huasen, R. Srikant, Xin Liu, Chong Jiang. 2015. Algorithms with logarithmic or sublinear regret for constrained contextual bandits. C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, R. Garnett, eds., *Advances in Neural Information Processing Systems* 28. Curran Associates, Inc., 433–441.

- Wu, Wei Biao. 2005. Nonlinear system theory: Another look at dependence. *Proceedings of the National Academy of Sciences* 102(40) 14150–14154.
- Yu, Bin. 1994. Rates of convergence for empirical processes of stationary mixing sequences. *Ann. Probab.* 22(1) 94–116.
- Zhang, Danna, Wei Biao Wu. 2017. Gaussian approximation for high dimensional time series. *Ann. Statist.* 45(5) 1895–1919.
- Zhang, Xianyang, Guang Cheng. 2014. Bootstrapping high dimensional time series. *ArXiv e-print 1406.1037* .
- Zhang, Xianyang, Guang Cheng. 2018. Gaussian approximation for high dimensional vector under physical dependence. *Bernoulli* 24(4A) 2640–2675.