

# The Big Data Newsvendor: Practical Insights from Machine Learning

Gah-Yi Ban\*

Management Science & Operations, London Business School, Regent's Park, London, NW1 4SA, United Kingdom.  
gban@london.edu

Cynthia Rudin

Department of Computer Science, Duke University, Durham, NC 27708, United States of America  
cynthia@cs.duke.edu

We investigate the data-driven newsvendor problem when one has  $n$  observations of  $p$  features related to the demand as well as historical demand data. Rather than a two-step process of first estimating a demand distribution then optimizing for the optimal order quantity, we propose solving the “Big Data” newsvendor problem via single step machine learning algorithms. Specifically, we propose algorithms based on the Empirical Risk Minimization (ERM) principle, with and without regularization, and an algorithm based on Kernel-weights Optimization (KO). The ERM approaches, equivalent to high-dimensional quantile regression, can be solved by convex optimization problems and the KO approach by a sorting algorithm. We analytically justify the use of features by showing that their omission yields inconsistent decisions. We then derive finite-sample performance bounds on the out-of-sample costs of the feature-based algorithms, which quantify the effects of dimensionality and cost parameters. Our bounds, based on algorithmic stability theory, generalize known analyses for the newsvendor problem without feature information. Finally, we apply the feature-based algorithms for nurse staffing in a hospital emergency room using a data set from a large UK teaching hospital and find that (i) the best ERM and KO algorithms beat the best practice benchmark by 23% and 24% respectively in the out-of-sample cost, and (ii) the best KO algorithm is faster than the best ERM algorithm by three orders of magnitude and the best practice benchmark by two orders of magnitude.

*Key words:* big data, newsvendor, machine learning, Sample Average Approximation, statistical learning theory, quantile regression

*History:* First version: February 6, 2014; revisions submitted February 1, 2015, August 7, 2016 and November 2, 2017. Accepted for publication in Operations Research on March 2, 2018.

---

## 1. Introduction

We investigate the newsvendor problem when one has access to  $n$  past demand observations as well as a potentially large number,  $p$ , of *features* about the demand. By features we mean exogenous variables (also known as covariates, attributes and explanatory variables) that are predictors of the demand and are available to the decision-maker (hereafter, “DM”) before the ordering occurs. While inventory models to date have typically been constructed with demand as the stochastic primitive,

\* Formerly: Gah-Yi Vahn.

in the world of Big Data, the DM has access to a potentially large amount of relevant information such as customer demographics, weather forecast, seasonality (e.g. day of the week, month of the year and season), economic indicators (e.g. the consumer price index), as well as past demands to inform her ordering decisions.

In this paper, we propose solving the “Big Data” newsvendor problem via distribution-free, one-step machine learning algorithms that handle high-dimensional feature data, and derive finite-sample performance bounds on their out-of-sample costs. The one-step algorithms contrast with the approach of solving the Big Data newsvendor problem via a two-step process, of first estimating a (feature-dependent) demand distribution then optimizing for the optimal order quantity. Such two-step processes can be problematic because demand model specification is difficult in high-dimensions, and errors in the first step will amplify in the optimization step.

In this setting, the paper is structured to answer the following four questions. (Q1) How should the DM use a feature-demand data set to solve the newsvendor problem? (Q2) What is the value of incorporating features in newsvendor decision-making in the first place? (Q3) What theoretical guarantees does the DM using such data have, and how do these scale with the various problem parameters? (Q4) How do newsvendor decisions based on the feature-demand data set compare to other benchmarks in practice?

We address (Q1) in Sec. 2, where we propose one-step approaches to finding the optimal order quantity with a data set of both demand and related feature observations. One approach is based on the machine learning principle of Empirical Risk Minimization (ERM) (Vapnik 1998), with and without regularization, and the other, which we call Kernel-weights Optimization (KO), is inspired by Nadaraya-Watson’s kernel regression method (Nadaraya 1964, Watson 1964). The ERM approach, equivalent to high-dimensional quantile regression, is a linear programming (LP) algorithm without regularization, and a mixed-integer program (MIP), an LP or a quadratic program (QP) with  $\ell_0, \ell_1$  and  $\ell_2$  regularizations respectively. Under the KO approach, the optimal data-driven order quantity can be found by a simple sorting algorithm.

We justify the use of features by answering (Q2) in Sec. 3, where we analytically quantify the value of features by comparing against decisions made without any features (the “Sample Average Approximation (SAA)” method). This is necessary as most data-driven inventory works to date do not consider features. We consider two demand models for the comparison — a two-population model and the linear demand model. For both models, we show that the no-feature SAA decision does not converge to the true optimal decision, whereas the feature-based ERM decision does. In other words, the SAA decision can have a constant bias (i.e.  $O(1)$  error) regardless of how many observations  $n$  the DM has, whereas any finite-sample bias of the feature-based decision shrinks to

zero as  $n$  tends to infinity. Accordingly, the SAA decision can have a higher newsvendor cost than the ERM decision. We quantify the additional cost incurred by biased decisions in Theorem 4.

We address (Q3) in Sec. 4. In practice, the DM may never make the truly optimal decision with finite amount of data, even if s/he uses as much relevant feature information as possible. To understand the tradeoffs of not having the full distributional information of the demand, we derive theoretical performance bounds for the DM who uses the algorithms proposed in Sec. 2. Our bounds characterize how the in-sample cost of the in-sample decision (which the DM can calculate, as opposed to the expected cost of the in-sample decision) deviates from the true expected cost of the true optimal decision in terms of the various parameters of the problem, in particular in terms of  $p$  and  $n$ . The bounds show how the out-of-sample cost of the in-sample decision (the “generalization error”) deviates from its in-sample cost by a complexity term that scales gracefully as  $1/\sqrt{n}$  for  $n$ , and as  $\sqrt{\log(1/\delta)}$  for  $\delta$ , where  $1 - \delta$  is the probabilistic accuracy of our bound, under the minimal assumption that the demand data is independent and identically distributed (iid) in the high-dimensional feature space. The bounds also show how the finite-sample bias from the true optimal decision scales as  $n^{-1/(2+p/2)}\sqrt{\log n}$ , under the additional assumption of linear demand model.

From a practical perspective, our bounds explicitly show the trade-offs between the generalization error (which measures the degree of in-sample overfitting) and the finite-sample bias and how they depend on the size of the data set, the newsvendor cost parameters and any free parameters (controls) in the algorithms. While some past papers in operations management have incorporated specific features in inventory models, none have analyzed the out-of-sample performance (cost) with high-dimensional data. Our work also contrasts with past works in quantile regression by using algorithm-specific stability theory of Bousquet and Elisseeff (2002) to analyze the generalization error, as opposed to Vapnik-Chervonenkis (VC) theory (Vapnik 1998), which results in bounds with tight constants dependent only on the newsvendor cost parameters, as opposed to ones dependent on uniform complexity measures such as covering numbers, VC dimensions and Rademacher averages, which are often difficult to compute and interpret in practice. Detailed discussions of the literature can be found below in Section 1.1.

We address (Q4) in Sec. 5, where we evaluate our algorithms against other known benchmarks through an extensive empirical investigation. Specifically, we apply our algorithms and other methods to a nurse staffing problem in a hospital emergency room. The best result using the ERM approach was with  $\ell_1$  regularization, with a cost improvement of 23% [a saving of £44,219 per annum (p.a.)] relative to the best practice benchmark (Sample Average Approximation (SAA) clustered by day of the week), and the best result using the KO approach had a cost improvement of 24% (a saving of £46,555 p.a.) relative to the same benchmark. Both results were statistically significant

at the 5% level. The best KO method, solved using the sorting procedure described in Sec. 2, was also very computationally efficient, taking just 0.05 seconds to compute the optimal staffing level for the next period, which is three orders of magnitude faster than the best ERM method and two orders of magnitude faster than the best practice benchmark and other benchmarks.

Finally, in Sec. 6, we conclude with a discussion of the practical take-aways, generalizable insights as well as limitations of our investigation, and thoughts on future directions for research.

### 1.1. Literature Review

Our work contributes to the following areas of investigation in operations management and machine learning.

(i) *Data-driven inventory models.* Models in newsvendor/inventory management have long been constructed with the available (or, rather, the lack of) data in mind, with the stochastic nature of the demand modelled with various assumptions. The earliest papers (Arrow et al. 1958, Scarf 1959b) make the assumption that the demand distribution is fully known, and this has been relaxed in recent years. Overall, there have been three main approaches to modeling demand uncertainty in the literature. In the Bayesian approach, the demand distribution is assumed to be known up to unknown parameter(s), which are dynamically learned starting with prior assumptions (Scarf 1959a, Azoury 1985, Lovejoy 1990). In the minimax approach, the decision-maker (hereafter, DM) opts for the best robust decision among all demand distributions within a specified uncertainty set (Scarf et al. 1958, Gallego and Moon 1993, Chen et al. 2007, Perakis and Roels 2008). Finally, in the data-driven approach, the DM has access to samples of demand observations drawn from an unknown distribution. In this setting, Burnetas and Smith (2000), Huh and Rusmevichientong (2009) and Kunnumkal and Topaloglu (2008) propose stochastic gradient algorithms, Godfrey and Powell (2001) and Powell et al. (2004) consider the adaptive value estimation method, Levi et al. (2007) provide a sampling-based method based on solving a shadow problem to solve for the optimal ordering quantities, and Levi et al. (2015) improve upon the bounds of Levi et al. (2007) for the single-period, featureless newsvendor case.

Within this line of work, the distinguishing aspect of our paper is in the incorporation and analysis of a potentially large number of features directly in an inventory model. As far as we are aware, the works of Liyanage and Shanthikumar (2005), See and Sim (2010) and Hannah et al. (2010) are the only other preceding methodological papers in operations management to incorporate some form of feature information in the decision model. We compare them with our approaches in detail in Sec. 2.4. Most importantly, we derive performance bounds that quantify the effect of the feature dimension on the out-of-sample cost, which has no precedence in this line of literature.

Our setup and subsequent analysis also differ from the parametric modeling approaches of Feldman (1978), Lovejoy (1992), Song and Zipkin (1993), Gallego and Özer (2001), Lu et al. (2006), Iida and Zipkin (2006), where the demand is modeled as Markov-modulated processes (MMDP) with known state-dependent distributions, where the states capture various exogenous information such as economic indicators or advanced demand information. The key difference is that we make minimal assumptions about the underlying demand distribution (iid for the most part and later, for more specific finite-sample bias analysis, the linear demand model with unknown error distribution), whereas in all of the aforementioned works, the demand is modelled parametrically with known, state-dependent distributions (e.g. Normal or Poisson, where the mean is a function of the state), with the state evolving as a Markov process.

(ii) *Performance bound analysis of data-driven inventory decisions.* While there is no precedent for feature-dependent performance bounds in the inventory theory literature, Levi et al. (2007) and Levi et al. (2015) provide such bounds without features. Levi et al. (2007) studies the single-period and dynamic inventory problems with zero setup cost from a nonparametric perspective and provides sample average approximation (SAA)-type algorithms to solve them. They then provide probabilistic bounds on the minimal number of iid demand observations that are needed for the algorithms to be near-optimal. Levi et al. (2015) improves upon the bound for the single-period case. Our performance bounds are generalizations of the bounds of Levi et al. (2007) and Levi et al. (2015) to incorporate features in the DM's data set. We demonstrate how our bounds can retrieve the bounds of Levi et al. (2007) when  $p = 1$  in Appendix D; a similar result can be shown for Levi et al. (2015) as well.

(iii) *High-dimensional Quantile Regression.* Due to the equivalence of the newsvendor cost function with the loss function in quantile regression, our work can be classified as a study in high-dimensional quantile regression as well; albeit with the twist that we analyze (analytically and empirically) the cost of the estimated quantile as opposed to the estimated quantile itself. We make two contributions to this literature. First, the KO method is, to the best of our knowledge, a new nonparametric quantile regression method. Second, our out-of-sample performance analyses of both the ERM methods and the KO method, which uses algorithmic stability theory from machine learning (Bousquet and Elisseeff 2002), are new. Lastly, we extend the results of Chaudhuri et al. (1991) to the high-dimensional setting to derive bounds on the biases of the newsvendor algorithms under consideration. For references on quantile regression, we refer the readers to Koenker (2005) for a textbook reference, Takeuchi et al. (2006), Chernozhukov and Hansen (2008), Chernozhukov et al. (2010), Belloni and Chernozhukov (2011) and references therein for more recent works on high-dimensionality.

## 2. Solving the Newsvendor Problem with Feature Data

### 2.1. The Newsvendor Problem

A company sells perishable goods and needs to make an order before observing the uncertain demand. For repetitive sales, a sensible goal is to order a quantity that minimizes the total expected cost according to:

$$\min_{q \geq 0} EC(q) := \mathbb{E}[C(q; D)], \quad (1)$$

where  $q$  is the order quantity,  $D \in \mathcal{D}$  is the uncertain (random) future demand,

$$C(q; D) := b(D - q)^+ + h(q - D)^+ \quad (2)$$

is the random cost of order  $q$  and demand  $D$ , and  $b$  and  $h$  are respectively the unit backordering and holding costs. If the demand distribution,  $F$ , is known, one can show the optimal decision is given by the  $b/(b + h)$  quantile, that is:

$$q^* = \inf \left\{ y : F(y) \geq \frac{b}{b + h} \right\}. \quad (3)$$

### 2.2. The Data-Driven Newsvendor Problem

In practice, the decision maker does not know the true distribution. If one has access to historical demand observations  $\mathbf{d}(n) = [d_1, \dots, d_n]$ , but no other information, then a sensible approach is to substitute the true expectation with a sample average expectation and solve the resulting problem:

$$\min_{q \geq 0} \hat{R}(q; \mathbf{d}(n)) = \frac{1}{n} \sum_{i=1}^n [b(d_i - q)^+ + h(q - d_i)^+], \quad (\text{SAA})$$

where we use the  $\hat{\cdot}$  notation to emphasize quantities estimated from data. This approach is called the Sample Average Approximation (SAA) approach in stochastic optimization (for further details on the SAA approach in stochastic optimization, see [Shapiro et al. 2009](#)). One can show the optimal SAA decision is given by

$$\hat{q}_n = \inf \left\{ y : \hat{F}_n(y) \geq \frac{b}{b + h} \right\}, \quad (4)$$

where  $\hat{F}_n(\cdot)$  is the empirical cdf of the demand from the  $n$  observations. Note that if  $F$  is continuous, and we let  $r = b/(b + h)$ , then  $\hat{q}_n = d_{[nr]}$ , the  $[nr]$ -th largest demand observation.

### 2.3. The Feature-Based Newsvendor Problem

In practice, the demand depends on many observable *features* (equivalently, independent/explanatory variables, attributes or characteristics), such as seasonality (day, month, season), weather, location and economic indicators, which are available prior to making the order. In other words, the real newsvendor problem is optimize the *conditional* expected cost function:

$$\min_{q(\cdot) \in \mathcal{Q}, \{q: \mathcal{X} \rightarrow \mathbb{R}\}} \mathbb{E}[C(q(\mathbf{x}); D(\mathbf{x})) | \mathbf{x}], \quad (5)$$

where the decision is now a function that maps the feature space  $\mathcal{X} \subset \mathbb{R}^p$  to the reals and the expected cost that we minimize is now conditional on the feature vector  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^p$ .

The decision-maker intent on finding an optimal order quantity in this new setting has three issues to address. The first issue is in knowing what features the demand depends on, which prescribes what data to collect. As this is application-specific, we assume that the decision maker has already collected appropriate historical data  $S_n = [(\mathbf{x}_1, d_1), \dots, (\mathbf{x}_n, d_n)]$ . The data may be low-dimensional, where the number of features  $p$  is small compared to the number of observations  $n$ , or high dimensional, where the number of features is large compared to  $n$  (the analysis of identifying the low- and high-dimensional regimes later in Sec. 4). The second issue is how to solve the problem (5) in an efficient manner given the feature-demand data set. In this section, we propose two approaches to solving (5) — the ERM and KO approaches. Both approaches are direct, in that the decision-maker solves for the (in-sample) optimal order quantity in a single step. As such, our proposed algorithms are customized for the feature-based newsvendor problem, and are distinct from SAA (which are independent of features) and the separated estimation and optimization (SEO) approach, against which we compare in Sec. 5, along with other known benchmarks. The final concern is what performance guarantee is possible prior to observing the demand in the next period. We address this in Sec. 4.

Before we begin, we clarify that the setting of interest is one in which the DM observes the features  $\mathbf{x}_{n+1}$  before making the ordering decision at time  $n + 1$ .

**2.3.1. Empirical Risk Minimization Algorithms.** The ERM approach to solving the newsvendor problem with feature data is:

$$\min_{q(\cdot) \in \mathcal{Q}, \{q: \mathcal{X} \rightarrow \mathbb{R}\}} \hat{R}(q(\cdot); S_n) = \frac{1}{n} \sum_{i=1}^n [b(d_i - q(\mathbf{x}_i))^+ + h(q(\mathbf{x}_i) - d_i)^+], \quad (\text{NV-ERM})$$

where  $\hat{R}$  is called the *empirical risk* of function  $q$  with respect to the data set  $S_n$ .

To solve (NV-ERM), one needs to specify the function class  $\mathcal{Q}$ . The size or the complexity of  $\mathcal{Q}$  controls overfitting or underfitting: for instance, if  $\mathcal{Q}$  is too large, it will contain functions that fit the noise in the data, leading to overfitting. Let us consider linear decision rules of the form

$$\mathcal{Q} = \left\{ q: \mathcal{X} \rightarrow \mathbb{R} : q(\mathbf{x}) = \mathbf{q}'\mathbf{x} = \sum_{j=1}^p q^j x^j \right\},$$

where  $x^1 = 1$ , to allow for a feature-independent term (an intercept term). This is not restrictive, as one can easily accommodate nonlinear dependencies by considering nonlinear transformations of basic features. We might, for instance, consider polynomial transformations of the basic features, e.g.,  $[x_1, \dots, x_p, x_1^2, \dots, x_p^2, x_1x_2, \dots, x_{p-1}x_p]$ . Such transformations can be motivated from generative models of the demand (but do not need to be); for instance, assume:  $D = f(\mathbf{x}) + \varepsilon$ , where  $\mathbf{x}$  is a  $p$ -dimensional vector of features. If we also assume that  $f(\cdot)$  is analytic, we can express the demand function by its Taylor expansion:

$$\begin{aligned} D &\approx f(\mathbf{0}) + \partial f(\mathbf{0})' \mathbf{x} + \mathbf{x}' [D^2 f(\mathbf{0})] \mathbf{x} + \dots + \varepsilon \\ &= f(\mathbf{0}) + \sum_{i=1}^p \partial f_i(\mathbf{0}) x_i + \sum_{i=1}^p \sum_{j=1}^p [D^2 f(\mathbf{0})]_{ij} x_i x_j + \dots + \varepsilon, \end{aligned} \quad (6)$$

which means that the demand function of a basic feature vector  $\mathbf{x}$  can be approximated by a linear demand model with a much larger feature space. For example, the second-order Taylor approximation of the demand model can be considered to be a linear demand model with the  $(p + p^2)$  features mentioned earlier:  $[x_1, \dots, x_p, x_1^2, \dots, x_p^2, x_1x_2, \dots, x_{p-1}x_p]$ . Regardless of the motivation for the transformations of basic features, we can choose them to be arbitrarily complex; hence our choice of decision functions that depend linearly on the feature vector is not particularly restrictive. The choice of  $\mathcal{Q}$  can be made more or less complex depending on which transformations are included. For a discussion on how piecewise linear decisions can be considered by transforming the features, see Sec. 2.4.2. This comes at the cost of increasing the number of features, perhaps dramatically so, thereby increasing the likelihood of overfitting if there are not enough data. We thus propose ERM with regularization for large  $p$  (discussed further in Sec. 4). Although we consider linear decision functions for the rest of the paper, we show how one can consider nonlinear functional spaces for  $\mathcal{Q}$  in Appendix A via mapping the original features onto a higher dimensional reproducing kernel Hilbert space.

When  $p$  is relatively small, the DM can solve (NV-ERM) via the following linear program:

#### ERM Algorithm 1

$$\begin{aligned} \min_{q: q(\mathbf{x}) = \sum_{j=1}^p q^j x^j} \quad & \hat{R}(q(\cdot); S_n) = \frac{1}{n} \sum_{i=1}^n [b(d_i - q(\mathbf{x}_i))^+ + h(q(\mathbf{x}_i) - d_i)^+] \\ \equiv \quad & \min_{\mathbf{q}=[q^1, \dots, q^p]} \frac{1}{n} \sum_{i=1}^n (bu_i + ho_i) \\ \text{s.t. } \forall i = 1, \dots, n: \quad & u_i \geq d_i - q^1 - \sum_{j=2}^p q^j x_i^j \\ & o_i \geq q^1 + \sum_{j=2}^p q^j x_i^j - d_i \\ & u_i, o_i \geq 0, \end{aligned} \quad (\text{NV-ERM1})$$



where the indicator variables  $u_i$  and  $o_i$  represent, respectively, underage and overage costs in period  $i$ . This is an LP with a  $p + 2n$ -dimensional decision vector and  $4n$  constraints. We will see in Sec. 4 that while (NV-ERM1) yields decisions that are algorithmically stable, the performance guarantee relative to the true optimal decision is loose when  $p$  is large (relative to  $n$ ). Thus, in the case of high dimensional data, one can solve the LP (NV-ERM1) by selecting a subset of the most relevant features according to some feature-selection criterion, for example via cross validation or via model selection criteria such as the Akaike Information Criterion (Akaike 1974) or Bayesian Information Criteria (Schwarz 1978). Alternatively, one can automate feature selection by solving the following *regularized* version of (NV-ERM1):

**ERM Algorithm 2 (with regularization)**

$$\begin{aligned}
\min_{q: q(\mathbf{x}) = \sum_{j=1}^p q^j x^j} \quad & \hat{R}(q(\cdot); S_n) + \lambda \|q\|_2^2 = \frac{1}{n} \sum_{i=1}^n [b(d_i - q(\mathbf{x}_i))^+ + h(q(\mathbf{x}_i) - d_i)^+] + \lambda \|\mathbf{q}\|_k^2 \\
\equiv \quad & \min_{\mathbf{q}=[q^1, \dots, q^p]} \frac{1}{n} \sum_{i=1}^n (bu_i + ho_i) \\
s.t. \quad & \forall i = 1, \dots, n: \\
& u_i \geq d_i - q^1 - \sum_{j=2}^p q^j x_i^j \\
& o_i \geq q^1 + \sum_{j=2}^p q^j x_i^j - d_i \\
& u_i, o_i \geq 0,
\end{aligned} \tag{NV-ERM2}$$

where  $\lambda > 0$  is the regularization parameter and  $\|\mathbf{q}\|_k$  denotes the  $\ell_k$ -norm of the vector  $\mathbf{q} = [q^1, \dots, q^p]$ . If we regularize by the  $\ell_2$  norm, the problem becomes a quadratic program (QP) and can be solved efficiently using widely available conic programming solvers. If we believe that the number of features involved in predicting the demand is very small, we can choose to regularize by the  $\ell_0$  semi-norm or the  $\ell_1$  norm to encourage sparsity in the coefficient vector. The resulting problem then becomes, respectively, a mixed-integer program (MIP) or an LP. Note regularization by the  $\ell_k$ -norm is widely used across engineering, statistics and computer science to handle overfitting.

Let us consider variations. We may want a set of coefficients to be either all present or all absent, for instance if they fall into the same category (e.g., all are weather-related features). We can accommodate this with a regularization term  $\sum_{g=1}^G \|q_{\mathcal{I}_g}\|_2$ , with  $\mathcal{I}_g$  being the indicator of group  $g$ . This regularization term is an intermediate between  $\ell_1$  and  $\ell_2$  regularization, where sparsity at the group level is encouraged by the sum ( $\ell_1$  norm) over groups. We will see in Sec. 4 that regularization leads to stable decisions with good finite-sample performance guarantees.

**2.3.2. Kernel Optimization (KO) Method.** Here we introduce an alternative approach that can take features into account. We call this approach the Kernel-weights Optimization (KO) method because it is based on Nadaraya-Watson kernel regression (Nadaraya 1964, Watson 1964).

One of the goals of nonparametric regression is to estimate the expectation of a dependent variable (e.g. demand) conditional on independent variables taking on a particular value. That is, given past data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  one wants to estimate

$$m(\mathbf{x}_{n+1}) = \mathbb{E}[Y|\mathbf{x}_{n+1}],$$

where  $Y \in \mathbb{R}$  is the dependent variable and  $\mathbf{x}_{n+1} \in \mathbb{R}^p$  is a vector of new independent variables. In 1964, Nadaraya and Watson proposed to estimate this quantity by the locally weighted average

$$m_h(\mathbf{x}_{n+1}) = \frac{\sum_{i=1}^n K_w(\mathbf{x}_{n+1} - \mathbf{x}_i) y_i}{\sum_{i=1}^n K_w(\mathbf{x}_{n+1} - \mathbf{x}_i)},$$

where  $K_w(\cdot)$  is a kernel function with bandwidth  $w$ . Typical examples of the kernel function include the uniform kernel

$$K(\mathbf{u}) = \frac{1}{2} \mathbb{I}(\|\mathbf{u}\|_2 \leq 1)$$

and the Gaussian kernel

$$K(\mathbf{u}) = \frac{1}{\sqrt{2\pi}} \exp^{-\|\mathbf{u}\|_2^2/2},$$

with  $K_w(\cdot) := K(\cdot/w)/w$ . Now for an order quantity  $q$ , the feature-dependent newsvendor expected cost after observing features  $\mathbf{x}_{n+1}$  is given by

$$\mathbb{E}[C(q; D)|\mathbf{x}_{n+1}], \tag{7}$$

which depends (implicitly) on the demand distribution at  $\mathbf{x}_{n+1}$ . Thus if we consider the newsvendor cost to be the dependent variable, we can estimate (7) by the Nadaraya-Watson estimator

$$\frac{\sum_{i=1}^n K_w(\mathbf{x}_{n+1} - \mathbf{x}_i) C(q, d_i)}{\sum_{i=1}^n K_w(\mathbf{x}_{n+1} - \mathbf{x}_i)}.$$

This gives rise to a new approach to feature-data-driven newsvendor, which we call the Kernel Optimization (KO) Method.

$$\min_{q \geq 0} \tilde{R}(q; S_n, \mathbf{x}_{n+1}) = \min_{q \geq 0} \frac{\sum_{i=1}^n K_w(\mathbf{x}_{n+1} - \mathbf{x}_i) C(q, d_i)}{\sum_{i=1}^n K_w(\mathbf{x}_{n+1} - \mathbf{x}_i)}. \tag{NV-KO}$$

Note that there are no edge effects in the objective estimate if the kernel is smooth, which is the case for the Gaussian kernel. Notice that the optimization is over the non-negative reals, and the optimal decision implicitly depends on  $\mathbf{x}_{n+1}$ . (NV-KO) is a one-dimensional piecewise linear optimization problem, and we can find its solution according to the following proposition.

PROPOSITION 1. *The optimal feature-based newsvendor decision  $\hat{q}_n^\kappa$  obtained by solving (NV-KO) is given by*

$$\hat{q}_n^\kappa = \hat{q}_n^\kappa(\mathbf{x}_{n+1}) = \inf \left\{ q : \frac{\sum_{i=1}^n \kappa_i \mathbb{I}(d_i \leq q)}{\sum_{i=1}^n \kappa_i} \geq \frac{b}{b+h} \right\}, \quad (8)$$

where for simplicity we introduce  $\kappa_i = K_w(\mathbf{x}_{n+1} - \mathbf{x}_i)$ . In other words, we can find  $\hat{q}_n^\kappa$  by ranking the past demand in increasing order, and choosing the smallest value at which the inequality in (8) is satisfied.

Notice that the left hand side (lhs) of the inequality in (8) is similar to the empirical cdf of the demand, except that each past demand observation  $d_i$  is re-weighted by the distance of its corresponding feature  $\mathbf{x}_i$  to the current feature  $\mathbf{x}_{n+1}$ .

## 2.4. Comparison with other inventory papers that incorporate features

In the above, we had identified three systematic approaches to incorporating features for the newsvendor problem. However, incorporating exogenous information in inventory decision-making is not entirely new. In what follows, we compare and contrast the algorithms introduced thus far to past works that incorporate exogenous information in inventory decision-making.

**2.4.1. Comparison with Liyanage and Shanthikumar (2005).** Our first comparison is with operational statistics (OS), which was first introduced by Liyanage and Shanthikumar (2005). The idea behind OS is to integrate parameter estimation and optimization rather than separate them. Let us illustrate how OS works by an example similar to the one used in Liyanage and Shanthikumar (2005).

Suppose the true demand has an exponential distribution, i.e.  $D \sim \exp(1/\theta)$ , and that the decision maker has access to  $d_1, \dots, d_n$  observations of past data. Then with straightforward calculations, one can show first estimating then optimizing (“Separated Estimation and Optimization”, hereafter SEO) leads to the decision

$$\hat{q}_{SEO} = \log \left( \frac{b+h}{b} \right) \bar{d}_n,$$

where  $\bar{d}_n$  is the sample average of the demand. Now consider instead the decision

$$\hat{q}_{OS}^1(\alpha) = \alpha \bar{d}_n \quad (9)$$

parameterized by a constant  $\alpha > 0$ . The OS approach then picks  $\alpha$  by the following optimization:

$$\min_{\alpha \geq 0} \mathbb{E}_\theta [C(\hat{q}_{OS}^1(\alpha); D)]. \quad (10)$$

As  $\alpha = \log((b+h)/b)$  is a feasible solution of (10), this guarantees the OS decision to yield a true expected cost that is bounded above by the true expected cost of the SEO decision. In other words, by construction we have

$$\mathbb{E}_\theta[C(\hat{q}_{OS}^1(\alpha^*); D)] \leq \mathbb{E}_\theta[C(\hat{q}_{SEO}; D)], \quad (11)$$

where  $\alpha^*$  is the optimal parameter in (10). With some computations, one can show

$$\alpha^* = \left\lceil \left( \frac{b+h}{h} \right)^{1/n+1} - 1 \right\rceil n.$$

Liyanage and Shanthikumar (2005) also shows that one can also improve upon the SAA optimal decision in terms of the true expected cost by considering the decision

$$\hat{q}_{OS}^2(\alpha, \beta) = d_{\lceil \beta - 1 \rceil} + \alpha(d_{\lceil \beta \rceil} - d_{\lceil \beta - 1 \rceil}), \quad (12)$$

where  $\beta \in \{1, \dots, n\}$  and  $\alpha \geq 0$  are parameters to be chosen via

$$\min_{\alpha \geq 0, \beta \in \{1, \dots, n\}} \mathbb{E}_\theta[C(\hat{q}_{OS}^2(\alpha, \beta); D)]. \quad (13)$$

As the above example illustrates, OS takes insight from the form of the decision derived by other methods (e.g. SEO and SAA) and constructively improves upon them in terms of the true expected cost simply by considering a decision that is a *function* of past demand data rather than a scalar quantity. In the parlance of our feature-based approach, the OS method is essentially considering meaningful statistics of past demand data as *features*. However, there is an important difference between the OS approach and ours, and this is in the way the unknown coefficients (parameters) of the decision function are chosen. Under our decision-making paradigm, one would simply input the sample average of past demand and differences of order statistics of past demand as features and choose the coefficients that minimize the *in-sample average cost*. In contrast, OS is based on the premise that one knows the distributional family the demand belongs to, and thus is able to compute the coefficients that minimize the *true expected cost*. That one knows the true distributional family is not a weak assumption, however the insights from OS analysis are valuable. In Sec. 5, we will consider solving (NV-ERM1) and (NV-ERM2) both without and with OS-inspired features, to evaluate their practical benefit in terms of the out-of-sample cost.

**2.4.2. Comparison with See and Sim (2010).** See and Sim (2010) investigates the multi-period inventory management problem in the presence of features such as “market outlook, oil prices, trend, seasonality, cyclic variation” that a product demand can depend on. Specifically, they model the demand at time  $t$  as

$$d_t(\tilde{z}) = d_t^0 + \sum_{k=1}^{N_t} d_t^k \tilde{z}_k,$$

where  $\tilde{z} = [\tilde{z}_1, \dots, \tilde{z}_{N_t}]$  represents random features, and  $d_t^k$ ,  $k = 0, \dots, N_t$  are the coefficients. They then make a number of assumptions on the random features (zero mean, positive definite covariance matrix, bounded support set which is second-order conic representable, etc.), all of which are assumed to be known to the DM, and solve an approximation of the robust problem by considering linear decision rules (linear as a function of the features) as well as piecewise linear decision rules. See and Sim (2010) demonstrates that piecewise linear decision rules have the best performance.

While See and Sim (2010) certainly consider the presence of features in their problem setup, their work is distinctly different from ours on a number of fronts. First of all, See and Sim (2010) is about *robust* decision-making, as opposed to *data-driven* decision-making; and as such their theoretical results, while interesting, do not pertain to the data-driven questions we explore in this paper. Secondly, See and Sim (2010) assumes that the DM has access to a number of key statistics regarding the random features are known to the DM, and performance bounds for their approximate decisions are derived as a function of these statistics. In contrast, we do not make any assumptions about the data-generating process other than iid, and as such our performance analysis are independent of statistics that are presumed known. Thirdly, See and Sim (2010) do not study the effect of high-dimensionality (the “Big” in Big Data), which is the central theme of this paper.

Lastly, See and Sim (2010) makes an interesting observation that considering decisions that are piecewise linear functions of the underlying features perform better than linear decision rules. However, in this paper we consider only linear decision rules because non-linear decisions can be transformed into linear decision rules with the addition of new features. We had explained how to enlarge the feature space when the true decision is an analytic function of the features in Section 2.3.1. Let us now illustrate how to enlarge the feature space to allow for piecewise linear decisions. For simplicity, consider a single feature  $x$ . We can then construct new features  $\tilde{x}_1 = \mathbb{I}(x \leq c_1)$ ,  $\tilde{x}_2 = \beta \mathbb{I}(c_1 < x \leq c_2)$  and  $\tilde{x}_3 = \mathbb{I}(x > c_3)$  based on the “basic” feature  $x$ , and some pre-specified constants  $c_1$  and  $c_2$ . Then the corresponding linear decision rule

$$\mathbf{q} = q^0 + q^1 \tilde{x}_1 + q^2 \tilde{x}_2 + q^3 \tilde{x}_3$$

has the same piecewise linear structure as the decisions considered in See and Sim (2010). By incorporating a large number of breakage points, any piecewise linear decision can be approximated by linear decision rules arbitrarily well, although perhaps with the addition of a large number of “new” features corresponding to each breakage point. Thus we solely focus on the feature-based newsvendor problem with decisions that are linear functions of the feature vector, where the feature dimension may be large.

**2.4.3. Comparison with Hannah et al. (2010).** Hannah et al. (2010) consider the single-period stochastic optimization problem

$$\min_{x \in \mathcal{X}} \mathbb{E}_Z[F(x, Z)|S = s]$$

where  $x \in \mathcal{X}$  is the decision variable,  $Z$  is a state-dependent random variable, and  $S = s$  is the current state of the world. They propose solving the above problem by the weighted empirical stochastic optimization problem

$$\min_{x \in \mathcal{X}} \sum_{i=1}^n w_n(s, S_i) F(x, Z(S_i)),$$

where  $w_n(s, S_i)$  are the weights given to the past data  $(S_i, Z(S_i))_{i=1}^n$  determined by the Nadaraya-Watson-based kernel estimator as in the (KO) method or by a complex Dirichlet process mixture model.

This is similar to our feature-based newsvendor setup, however there are a few key differences. Firstly, the model studied by Hannah et al. (2010) requires discrete state variables, whereas we consider both discrete and continuous feature variables. Secondly, Hannah et al. (2010) do not study the effect of high-dimensionality (the “Big” in Big Data), which is the central theme of this paper (their numerical example, which also considers the newsvendor problem, has only two states). Finally, Hannah et al. (2010) report numerical studies of computing the in-sample decisions, whereas we provide both theoretical performance guarantees and extensive empirical computation of out-of-sample performance of the in-sample decisions.

### 3. Value of Feature Information

In this section, we quantify the value of incorporating features in newsvendor decision-making, by comparing the (NV-ERM1) decision against the SAA decision, which is made with only past demand data. We consider two demand models: a two-population demand model, and the linear demand model. In both cases, the SAA decision yields inconsistent decisions (i.e. the in-sample decision does not converge to the true optimal, even when given an infinite amount of iid demand data), whereas the feature-based decision is consistent. We further quantify the implication on the expected cost, which is necessarily larger for decisions that are further away from the true optimal. All proofs can be found in Appendix B.

#### 3.1. Motivating Example I: Two Population Model

Consider the following demand model:

$$D = D_0(1 - x) + D_1x, \tag{14}$$

where  $D_0$  and  $D_1$  are non-negative continuous random variables such that the corresponding critical newsvendor fractiles  $q_0^*$  and  $q_1^*$  follow  $q_0^* < q_1^*$ , and  $x \in \{0, 1\}$  is a binary feature (e.g. 0 for weekday and 1 for weekend, or 0 for male and 1 for female). Let  $p_0$  be the proportion of time  $x = 0$ . We have  $n$  historical observations:  $[(x_1, d_1), \dots, (x_n, d_n)]$ , of which  $n_0 = np_0$  are when  $x = 0$  and  $n_1 = n - n_0$  are when  $x = 1$  (assume rounding effects are negligible). Note the observations  $d_k$  can be decomposed into:  $\{d_k | x_k = 0\} = d_k^0$  and  $\{d_k | x_k = 1\} = d_k^1$ . Also let  $r = b/(b+h)$  to simplify notations. Let  $F_0$  and  $F_1$  denote the cumulative distribution functions (cdfs),  $F_0^{-1}$  and  $F_1^{-1}$  denote the inverse cdfs, and  $f_0$  and  $f_1$  the probability density functions (pdfs) of  $D_0$  and  $D_1$  respectively.

We also assume the following.

**Condition (A).** Assume  $F_0$  and  $F_1$  are twice differentiable (i.e.  $f_0$  and  $f_1$  are differentiable) and that there exists a  $0 < \gamma < 2$  such that

$$\sup_{0 < y < 1} y(1-y) \frac{|J_i(y)|}{f(F_i^{-1}(y))} \leq \gamma,$$

where  $J_i(\cdot)$  is the *score function* of distribution  $F_i$  defined by

$$J_i(y) = \frac{-f'_i(F_i^{-1}(y))}{f_i(F_i^{-1}(y))} = -\frac{d}{dy} \ln f_i(F_i^{-1}(y)).$$

Condition A is satisfied by many standard distributions such as uniform, exponential, logistic, normal, and log normal, for  $\gamma$  between 0 and  $\sim 1.24$ . The critical ratios for the uniform, exponential and logistic distributions can be computed straight-forwardly; for the normal distribution it is easier to compute the critical ratio by using the following equivalent formulation for the critical ratio:

$$\sup_{x \in \text{dom}(D)} F(x)(1-F(x)) \frac{|f'(x)|}{f(x)^2}.$$

In Table 1, we display some standard distributions that satisfy the requirement of Condition A. For more details see [Parzen \(1979\)](#).

Distribution	$f(F^{-1}(y))$	$J(y)$	Is $\sup_{0 < y < 1} y(1-y) \frac{ J(y) }{f(F^{-1}(y))} < 2$ ?
Uniform	1	0	Yes, LHS = 0
Exponential	$1-y$	1	Yes, LHS = 1
Logistic	$y(1-y)$	$2y-1$	Yes, LHS = 1
Normal	$\frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2} \Phi^{-1}(y) ^2\}$	$\Phi^{-1}(y)$	Yes, LHS = 1
Lognormal	$\phi(\Phi^{-1}(y)) \exp\{-\Phi^{-1}(y)\}$	$\exp\{-\Phi^{-1}(y)\}(\Phi^{-1}(y)+1)$	Yes, LHS $\lesssim 1.24$

**Table 1** Some standard distributions that satisfy the requirement of Condition A. The standard normal cdf and pdf are denoted as  $\Phi(\cdot)$  and  $\phi(\cdot)$  respectively.

Below we derive the optimal in-sample solution given by (NV-ERM1).

LEMMA 1 (**Optimal ordering decision of (NV-ERM1)**). Let  $\hat{F}_i$  denote the empirical cdf of  $D|x=i$  with  $n_i$  iid observations for  $i=0,1$ . Then the optimal decision that solves (NV-ERM1) is given by

$$\begin{aligned}\hat{q}_n^0 &= \inf \left\{ q : \hat{F}_0(q) \geq \frac{b}{b+h} \right\} = d_{(\lceil n_0 r \rceil)}^0, \text{ if } x_{n+1} = 0 \\ \hat{q}_n^1 &= \inf \left\{ q : \hat{F}_1(q) \geq \frac{b}{b+h} \right\} = d_{(\lceil n_1 r \rceil)}^1, \text{ if } x_{n+1} = 1.\end{aligned}$$

Put simply,  $\hat{q}_n^0$  solves the SAA problem for the subsample of data corresponding to  $x=0$  and  $\hat{q}_n^0 + \hat{q}_n^1$  solves the SAA problem for the subsample of data corresponding to  $x=1$ .

THEOREM 1 (**Finite-sample bias and asymptotic optimality of (NV-ERM1)**). We can show

$$|\mathbb{E}[\hat{q}_n^i] - F_0^{-1}(r)| \leq O\left(\frac{\log n_i}{n_i}\right), \quad i=0,1$$

i.e. the finite-sample decision of the feature-based decision is biased by at most  $O(\log n_i/n_i)$ ,  $i=1,2$ , and

$$\lim_{n \rightarrow \infty} \hat{q}_n^i \stackrel{a.s.}{=} F_0^{-1}(r) =: q_i^*, \quad i=0,1$$

i.e. the feature-based decision is asymptotically optimal, correctly identifying the case when  $x=0$  or 1 as the number of observations goes to infinity.

LEMMA 2 (**Optimal SAA ordering decision**). Let  $F^{mix}$  denote the cdf of the mixture distribution  $D^{mix} = p_0 D_0 + (1-p_0) D_1$  and  $\hat{F}_n^{mix}$  its empirical counterpart with  $n$  observations. Then the optimal SAA decision is given by

$$\hat{q}_n^{SAA} = \inf \left\{ q : \hat{F}_n^{mix}(q) \geq \frac{b}{b+h} \right\} = d_{(\lceil nr \rceil)}.$$

THEOREM 2 (**Finite-sample bias and asymptotic (sub)-optimality of SAA**). The finite-sample bias of the SAA decision is given by

$$|\mathbb{E}[\hat{q}_n^{SAA}] - (F^{mix})^{-1}(r)| \leq O\left(\frac{\log n}{n}\right), \quad (15)$$

where  $(F^{mix})^{-1}$  is the inverse cdf of  $D^{mix}$ . Hence we also have

$$\begin{aligned}|\mathbb{E}[\hat{q}_n^{SAA} - \hat{q}_n^0]| &= |(F^{mix})^{-1}(r) - F_0^{-1}(r)| + O\left(\frac{\log n}{n}\right) = O(1) \\ |\mathbb{E}[\hat{q}_n^1 - \hat{q}_n^{SAA}]| &= |F_1^{-1}(r) - (F^{mix})^{-1}(r)| + O\left(\frac{\log n}{n}\right) = O(1).\end{aligned} \quad (16)$$

That is, on average, if  $x=0$  in the next decision period, the SAA decision orders too much and if  $x=1$  the SAA decision orders too little. In addition,

$$q_0^* < \lim_{n \rightarrow \infty} \hat{q}_n^{SAA} \stackrel{a.s.}{=} (F^{mix})^{-1}(r) < q_1^*, \quad (17)$$

hence the SAA decision is not asymptotically optimal (is inconsistent).



As a final point, we remark that these observations are analogous in spirit to the bias and inconsistency of regression coefficients when there are, in econometric parlance, correlated omitted variables in the model.

### 3.2. Motivating Example II: Linear Demand Model

Suppose the demand is given by the following linear model:

$$D|(\mathbf{X} = \mathbf{x}) = \beta^\top \mathbf{x} + \varepsilon, \quad (18)$$

where  $\varepsilon \sim F_\varepsilon$  is independent of the (random) feature vector  $\mathbf{X}$ , is continuous with probability density function  $f_\varepsilon(\cdot)$  which is bounded away from zero on the ordering domain  $[\underline{D}, \bar{D}]$  and has zero mean, and  $\mathbf{X}_1 = 1$  almost surely, to allow for a constant location term. In other words, the demand  $D$  depends linearly on the random features  $\mathbf{X} : \mathcal{X} \rightarrow \mathbb{R}^p$ , with some error. This is a widely-used and useful demand model that, apart from the fact that it can arbitrarily approximate nonlinear models as outlined in (6), it also subsumes times series models and the Martingale Model of Forecast Evolution (MMFE) of Heath and Jackson (1994) and Graves et al. (1986). For example, the autoregressive model of degree  $p$ ,  $AR(p)$  is modeled by

$$D_t = \alpha_0 + \alpha_1 D_{t-1} + \dots + \alpha_p D_{t-p} + \varepsilon_t,$$

where  $\varepsilon_t$  is a noise term with zero mean; this is clearly a linear demand model with features  $D_{t-p}, \dots, D_{t-1}$ . Also, the additive MMFE model for the demand at time  $t$  is given recursively by

$$D_t = D_{t-1} + \varepsilon_{t-1,t},$$

where  $\varepsilon_{t-1,t}$  is a mean zero normal random variable that captures forecast update at time  $t-1$  for demand at time  $t$ . Expanding the recursion, we get, at time 1:

$$D_t = D_0 + \sum_{i=0}^{t-1} \varepsilon_{i,i+1},$$

where  $D_0$  is known; and so the demand follows a linear model with  $\varepsilon_{i,i+1}$ ,  $i = 0, \dots, t-1$  as features.

A DM without the feature information only has access to past demand data:  $\mathcal{D} = \{d_1, \dots, d_n\}$ ; and a DM who has both past feature and demand data has the information:  $\mathcal{D}_\mathbf{x} = \{(\mathbf{x}_1, d_1), \dots, (\mathbf{x}_n, d_n)\}$ . Let  $\mathbf{x}_{n+1}$  denote the feature at time  $n+1$ , which is available to the DM. Then, the optimal order quantity is given by

$$q^*(\mathbf{x}_{n+1}, \mathbf{z}_{n+1}) = Q_\varepsilon \left( \frac{b}{b+h} \right) + \beta^\top \mathbf{x}_{n+1} \quad (19)$$

where

$$Q_\varepsilon\left(\frac{b}{b+h}\right) = \inf\{y : F_\varepsilon(y) \geq \frac{b}{b+h}\}$$

is the  $b/(b+1)$ -quantile of the distribution of  $\varepsilon$ .

The SAA solution is given by

$$\hat{q}^{SAA}(\mathbf{x}_{n+1}) = \inf\{y : \hat{F}_n^0(y) \geq \frac{b}{b+h}\} + \bar{d}_n \quad (20)$$

where  $\bar{d}_n$  is the sample average of the demand data, and  $\hat{F}_n^0$  is the empirical distribution of  $\{d_i - \bar{d}_n\}_{i=1}^n$ . Note the SAA decision is not dependent on any of the features, as the DM does not have access to any feature data. Thus the SAA decision is to order the same critical fractile quantity for the entire population at time  $n+1$ , regardless of what the population or the particular point in time may be. As a concrete example, this is like a national newspaper vendor who stocks the same number of newspapers at all shops, disregarding location features pertaining to the shop (e.g. customer demographics in the area) as well as relevant temporal features (e.g. holiday or not, weekday versus weekend, major political or sporting events) or features that pertain to both (e.g. historical demand for the shop).

The DM with features however orders the quantity

$$\hat{q}^{DM2}(\mathbf{x}_{n+1}) = \inf\{y : \hat{F}_n^2(y) \geq \frac{b}{b+h}\} + \sum_{k=2}^p \hat{q}^k x_{n+1}^k$$

and  $\hat{F}_n^2$  is the empirical distribution of  $\{d_i - \sum_{k=2}^p \hat{q}^k x_i^k\}$ , and  $\hat{q}^k$ ,  $k = 2, \dots, p$  are the solution coefficients to (NV-ERM1). Unlike the SAA decision, DM2's decision does depend on all relevant features. Continuing on with the national newspaper example, this corresponds to orders being different across stores as well as in time, taking into account such information as past sales and customer demographics at each store as well as temporal effects such as holidays/weekends and major public events.

**THEOREM 3 (Using no features leads to inconsistent decisions).** *Under the linear demand model of (18), given features  $\mathbf{X} = \tilde{\mathbf{x}}$ ,*

$$\begin{aligned} \hat{q}_n^{SAA}(\tilde{\mathbf{x}}) &\xrightarrow{a.s.} Q_\varepsilon\left(\frac{b}{b+h}\right) + \mathbb{E}_{\mathbf{X}}[\mathbb{E}_\varepsilon[D|\mathbf{X}]] \\ &= Q_\varepsilon\left(\frac{b}{b+h}\right) + \beta^\top \mathbb{E}[\mathbf{X}], \end{aligned}$$

and

$$\hat{q}_n^{DM2}(\tilde{\mathbf{x}}) \xrightarrow{a.s.} Q_\varepsilon\left(\frac{b}{b+h}\right) + \beta^\top \tilde{\mathbf{x}} = q^*(\tilde{\mathbf{x}}),$$

as  $n$  tends to infinity.

Considering the same national newspaper example above, we see that the SAA decision converges to the critical fractile of the dispersion  $\varepsilon$  plus the population-temporal average demand  $\mathbb{E}_{\mathbf{X}}[\mathbb{E}_{\varepsilon}[D|\mathbf{X}]]$ , whereas DM2's decision converges to the correct one,  $q^*(\tilde{\mathbf{x}})$ .

The results of Theorems 1–3 indicate that the no-feature SAA decisions are inconsistent, i.e., even with infinite amount of demand data the SAA decisions converge to quantities different from the true optimal. The natural question “is a decision that is further away from the true optimal necessarily worse in terms of the expected cost?” then follows. In other words, does the loss in the expected cost increase when the effect of the feature information increases? The answer is in the affirmative, which we detail below.

**THEOREM 4 (Expected Cost Difference).** *Let  $q^* = q^*(\tilde{\mathbf{x}})$  denote the true optimal newsvendor decision given feature  $\tilde{\mathbf{x}}$ , and  $\hat{q}$  some other decision not equal to  $q^*$ . Then the difference of the expected costs of the two decisions is given by*

$$\mathbb{EC}(\hat{q}; D) - \mathbb{EC}(q^*; D) = (b + h)\mathbb{E}[|\hat{q} - D|\mathbb{I}\{(\hat{q} \wedge q^*) \leq D \leq (\hat{q} \vee q^*)\}]. \quad (21)$$

Theorem 4 thus provides an exact formula for the expected cost difference of the sub-optimal decision  $\hat{q}$  from the true optimal decision  $q^*$ . We observe that the expected cost difference scales as the expectation of  $|\hat{q} - D|$  over the interval between the two decisions,  $\hat{q}$  and  $q^*$ . Thus what matters is the size of this interval and how the demand is distributed over it— the more concentrated the distribution over this interval, the larger the difference. While the exact quantity can only be computed with the knowledge of the demand distribution and the true optimal decision, we nevertheless arrive at the universal insight that the expected cost difference increases as  $\hat{q}$  deviates further away from  $q^*$ .

In particular, Theorem 4 implies that for the two population model (14), the more distinct the two population demands  $D_0$  and  $D_1$  in their critical fractiles, the worse the expected cost of the no-feature SAA decision to the true optimal solution. Likewise, for the linear demand model (18), the more idiosyncratic the feature information  $\beta^\top \tilde{\mathbf{x}}$  over the average  $\beta^\top \bar{\mathbf{X}}$ , the worse the expected cost of the no-feature SAA decision in comparison to the true optimal decision.

While Theorem 4 together with Theorems 1–3 justify the collection of features, a shortcoming is that the DM needs to know the demand distribution in order to quantify the gain in the expected cost due to a sub-optimal in-sample decision, which of course she does not know. In the next section, we characterize, with high probability bounds, the expected cost of the DM's in-sample decision using the information at hand.

#### 4. Bounds on the out-of-sample cost using in-sample information

In this section, we provide theoretical guarantees on the out-of-sample cost of the ordering decisions chosen by (NV-ERM1), (NV-ERM2) and (NV-KO). While (EC.4) from Theorem 4 states the exact expected cost difference between a proposed decision and the optimal decision, in practice one cannot compute the expectations because the demand distribution is unknown. The goal of this section is thus to provide performance bounds that are computable with the data at hand. We will see that the performance bound splits into two components; one that pertains to the bias arising from having a finite amount of data, commonly referred to as the *finite-sample bias*, and one that pertains to the variance of the in-sample decision, known as the generalization error.

The term “generalization” refers to the generalizability of the in-sample decision to out-of-sample data, and is a measure of the degree of over-fitting, as decisions with larger generalization errors are associated with greater over-fitting. Decisions that overfit can be misleading; for instance, if there is a large degree of freedom in the choice of the in-sample decision, then it may be possible to have zero in-sample cost leading the DM to think perfect ordering is possible. An astute DM who is aware of the perils of overfitting, however, would be cautious to make conclusions on the in-sample data alone, and the generalization error provides a description of how over-fitting arises for her decision.

The performance bounds derived in this section describe the mechanisms at play. They inform how the finite-sample bias and the generalization error scale with respect to the size of the data set at hand (i.e., with respect to  $p$  and  $n$ ), the problem parameters ( $b$  and  $h$ ) and any decision parameters (e.g. regularization parameter  $\lambda$  in (NV-ERM2)). As we will see, there is a tension between finite-sample bias and generalization error, which can be controlled by parameters such as the regularization parameter  $\lambda$  in (NV-ERM2) or the bandwidth parameter  $w$  in (NV-KO).

We start with some definitions. The *true risk* is the expected out-of-sample cost, where the expectation is taken over an unknown distribution over  $\mathcal{X} \times \mathcal{D}$ , where  $\mathcal{X} \subset \mathbb{R}^p$ . Specifically,

$$R_{true}(q) := \mathbb{E}_{D(\mathbf{x})}[C(q; D(\mathbf{x}))].$$

We are interested in minimizing this cost, but we cannot measure it as the distribution is unknown. The empirical risk is the average cost over the training sample:

$$\hat{R}(q; S_n) := \frac{1}{n} \sum_{i=1}^n C(q, d_i(\mathbf{x}_i)).$$

The empirical risk can be calculated, whereas the true risk cannot; the empirical risk alone, however, is an incomplete picture of the true risk. We must have some additional property of the algorithm to ensure that the method does not overfit. If the algorithm is stable, it is less likely to overfit, which we quantify in the results in this section. Specifically, we provide probabilistic upper bounds on the

true risk in terms of the empirical risk and the algorithmic stability of the method. Since we desire the true risk to be low, a combination of low empirical risk and sufficient stability ensures this.

The training set is, as before,  $S_n = \{z_1 = (\mathbf{x}_1, d_1), \dots, z_n = (\mathbf{x}_n, d_n)\}$ ,  $z \in \mathcal{Z}$ , and we also define the modified training set

$$S_n^{\setminus i} := \{z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n\},$$

which leaves one observation out.

A *learning algorithm* is a function  $A$  from  $\mathcal{Z}^n$  into  $\mathcal{Q} \subset \mathcal{D}^{\mathcal{X}}$ , where  $\mathcal{D}^{\mathcal{X}}$  denotes the set of all functions that map from  $\mathcal{X}$  to  $\mathcal{D}$ . A learning algorithm  $A$  maps the training set  $S_n$  onto a function  $A_{S_n} : \mathcal{X} \rightarrow \mathcal{D}$ . A learning algorithm  $A$  is *symmetric with respect to  $S_n$*  if for all permutations  $\pi : S_n \rightarrow S_n$  of the set  $S_n$ ,

$$A_{S_n} = A_{\pi(S_n)} = A_{\{\pi(z_1), \dots, \pi(z_n)\}}.$$

In other words, a symmetric learning algorithm does not depend on the order of the elements in the training set  $S_n$ . The *loss* of the decision rule  $q \in \mathcal{Q}$  with respect to a sample  $z = (\mathbf{x}, d)$  is defined as

$$\ell(q, z) := c(q, d(\mathbf{x})),$$

for some cost function  $c$ , which in our work is the newsvendor cost  $C$ .

We derive performance bounds on (NV-ERM1)–(NV-KO) under the following assumptions.

**Assumptions for Theorems 5–7.**

1. The feature vector  $\mathbf{X}$  is normalized ( $\mathbf{X}_1 = 1$  almost surely,  $\mathbf{X}_{[2:p]}$  has mean zero and standard deviation one) and that it lives in a closed unit ball:  $\|\mathbf{X}\|_2 \leq X_{\max} \sqrt{p}$ .
2. The demand follows the linear model (18) where the distribution of  $\varepsilon$ ,  $f_\varepsilon$ , is bounded away from zero on the domain  $[\underline{D}, \bar{D}]$  (otherwise unspecified).
3. All decision functions (policies) described are measurable, and  $\mathcal{Q}$  is a convex subset of a linear space.

Assumption 1 is for the feature vector  $\mathbf{X}$ ; we note the normalization assumption is to simplify the exposition and the results do not require that the DM knows the true mean or standard deviations of  $\mathbf{X}$ , only the size bound  $X_{\max}$ , the existence of which is realistic and not prohibitive. Assumption 2 details assumptions on the demand model, which is assumed to be linear for tractability, but we do not assume any distributional knowledge beyond its total range. Finally, Assumption 3 is a necessary requirement for sensible optimization over a function class when the demand is linear.

First, we state the performance bound on (NV-ERM1).

**THEOREM 5 (Out-of-sample performance of (NV-ERM1)).** *Denote the true optimal solution by  $q^* = q^*(\mathbf{x}_{n+1})$ , and the decision due to (NV-ERM1) by  $\hat{q} = \hat{q}(\mathbf{x}_{n+1})$ . Then with probability at least*

$1 - \delta$  over the random draw of the sample  $S_n$ , where each element of  $S_n$  is drawn i.i.d. from an unknown distribution on  $\mathcal{X} \times \mathcal{D}$ , and for all  $n \geq 3$ ,

$$|R_{true}(q^*) - \hat{R}_{in}(\hat{q}; S_n)| \leq (b \vee h) \bar{D} \left[ \frac{2(b \vee h)}{b \wedge h} \frac{p}{n} + \left( \frac{4(b \vee h)}{b \wedge h} p + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}} \right] + (b \vee h) K \frac{\sqrt{\log n}}{n^{1/(2+p/2)}},$$

where  $K = \sqrt{\frac{9(8+5p)}{(4+p)}} \frac{1}{(1-2^{-4/(4+p)})\lambda_2^*}$ , and  $\lambda_2^* = \min_{t \in [\underline{D}, \bar{D}]} f_\varepsilon(t)$ .

Theorem 5 is a statement about how close the in-sample cost of the in-sample decision,  $\hat{R}_{in}(\hat{q}; S_n)$  is to the expected cost of the true optimal decision,  $R_{true}(q^*)$ , in terms of quantities the DM knows.

The first term on the right hand side upper bound is the bound on the generalization error, which is the difference between the training error and test error for the in-sample decision. For fixed cost parameters  $b$  and  $h$ , we find that the generalization error scales as  $O(p/\sqrt{n})$ . Thus if the number of relevant features in the population model is small and not growing relative to the number of observations, then in-sample decisions generalize well to out-of-sample data; in other words overfitting should not be an issue. However, if  $p/\sqrt{n}$  is large, or growing (which happens when new observations are associated with new features), then overfitting will be an issue, and Theorem 5 suggests that (NV-ERM1) may not be a good algorithm in such a scenario. The dependence on the upper bound on the demand,  $\bar{D}$  is necessary so that the bound is not scale invariant. In other words, if the risks on the left hand side of the inequality changed units (e.g. from dollars per kilo of demand to dollars per ton), it would not make sense for the right hand side of the inequality to stay the same. Lastly, we note that the bound on the generalization error is tight in the sense that it comes from showing that the probability of large deviation of  $|\hat{R}_{true}(\hat{q}) - \hat{R}_{in}(\hat{q}; S_n)|$  decays exponentially fast in the number of observations  $n$ , which we establish through a property known as *uniform stability* of a learning algorithm, and because the constants in the bound are the smallest possible. For further details, we refer the reader to the proof in Appendix C and Bousquet and Elisseeff (2002). These are the best finite-sample bounds we know of for this problem. This is due to the fact that they are not uniform bounds, which require a complexity measure for the entire decision space (e.g. covering numbers, VC dimension or Rademacher complexity), rather algorithm-specific bounds that considers how the algorithm searches the decision space.

The second term on the right hand side upper bound is due to the finite-sample bias,  $\mathbb{E}|q^* - \hat{q}|$ . The only way a DM can reduce the finite-sample bias is by collecting more observations. The rate  $n^{-1/(2+p/2)}\sqrt{\log n}$  is optimal and cannot be improved upon without further assumptions on the demand model and/or the data generating process. For details, we refer the reader to the proof of Theorem 5 in Appendix C.

When  $p = 1$ , we are in the setup where the demand does not depend on any exogenous features (recalling that  $\mathbf{X}_1 = 1$  is the intercept term). This setting had been studied by Levi et al. (2007)

and our results are consistent with them in that their sampling bound, up to a constant factor, can be obtained from our bound. The details of this can be found in Appendix D.

We now state the performance bound on (NV-ERM2).

**THEOREM 6 (Out-of-sample performance of (NV-ERM2)).** *Denote the true optimal solution by  $q^* = q^*(\mathbf{x}_{n+1})$ , the decision due to (NV-ERM1) by  $\hat{q} = \hat{q}(\mathbf{x}_{n+1})$  and the decision due to (NV-ERM2) by  $\hat{q}_\lambda = \hat{q}_\lambda(\mathbf{x}_{n+1})$ . Then with probability at least  $1 - \delta$  over the random draw of the sample  $S_n$ , where each element of  $S_n$  is drawn i.i.d. from an unknown distribution on  $\mathcal{X} \times \mathcal{D}$ , and for all  $n \geq 3$ ,*

$$|R_{true}(q^*) - \hat{R}_{in}(\hat{q}_\lambda; S_n)| \leq (b \vee h) \bar{D} \left[ \frac{(b \vee h) X_{\max}^2 p}{n \lambda \bar{D}} + \left( \frac{2(b \vee h) X_{\max}^2 p}{\lambda \bar{D}} + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}} \right] \\ + (b \vee h) \mathbb{E}_{D|\mathbf{x}_{n+1}}[|\hat{q}_\lambda - \hat{q}|] + (b \vee h) K \frac{\sqrt{\log n}}{n^{1/(2+p/4)}},$$

where  $K = \sqrt{\frac{9(8+5p)}{(4+p)}} \frac{1}{(1-2^{-4/(4+p)})\lambda_2^*}$ , and  $\lambda_2^* = \min_{t \in [\underline{D}, \bar{D}]} f_\varepsilon(t)$ .

The performance bound of Theorem 6 has three components. The first term is a bound on the generalization error, which is of  $O(p/(\sqrt{n}\lambda))$ . Thus the amount of overfitting can be directly controlled by the amount of regularization imposed on the problem; the larger the  $\lambda$ , the smaller the generalization error. Choosing  $\lambda = O(1/p^2)$  retrieves the same error rate as for (NV-ERM1), so  $\lambda = O(1/p^2)$  is a good starting point for choosing  $\lambda$ . The bound thus provides a sense of the “right” scale for lambda, which is useful when you have to search for its best value in practice.

The second term, which does not appear in Theorem 5, is the bias of the in-sample decision due to regularization, in other words the bias due to having perturbed the optimization problem away from the true problem of interest. This term is larger for larger  $\lambda$ , and so there is an inherent trade-off between the generalization error and the regularization bias. Ultimately, however, regularization gives the DM an extra degree of control while being agnostic to which feature is important a priori, and in practice would work with the optimal value of  $\lambda$  that balances the generalization error-regularization bias tradeoff on a validation data set (see Sec. 5).

The third and the final term is the finite-sample bias. We note that while the regularization bias can be controlled by  $\lambda$ , the finite-sample bias can only be controlled by collecting more data.

Finally, we have the following result for (NV-KO).

**THEOREM 7 (Out-of-sample performance of (NV-KO)).** *Denote the true optimal solution by  $q^* = q^*(\mathbf{x}_{n+1})$ , the decision due to (NV-ERM1) by  $\hat{q} = \hat{q}(\mathbf{x}_{n+1})$  and the decision to (NV-KO) with the Gaussian kernel by  $\hat{q}^\kappa = \hat{q}^\kappa(\mathbf{x}_{n+1})$ . Then with probability at least  $1 - \delta$  over the random draw of*

the sample  $S_n$ , where each element of  $S_n$  is drawn i.i.d. from an unknown distribution on  $\mathcal{X} \times \mathcal{D}$ , and for all  $n \geq 3$ ,

$$|R_{true}(q^*) - \hat{R}_{in}(\hat{q}^\kappa; S_n)| \leq (b \vee h) \bar{D} \left[ \frac{2(b \vee h)}{b \wedge h} \frac{1}{1 + (n-1)r_w(p)} + \left( \frac{4(b \vee h)}{1/n + (1-1/n)r_w(p)} + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}} \right] \\ + (b \vee h) \mathbb{E}_{D|\mathbf{x}_{n+1}}[|\hat{q}^\kappa - \hat{q}|] + (b \vee h) K \frac{\sqrt{\log n}}{n^{1/(2+p/2)}},$$

where  $r_w(p) = \exp(-2X_{\max}^2 p/w^2)$ ,  $w$  the kernel bandwidth, and  $K = \sqrt{\frac{9(8+5p)}{(4+p)}} \frac{1}{(1-2-4/(4+p))\lambda_2^*}$ , and  $\lambda_2^* = \min_{t \in [\underline{D}, \bar{D}]} f_\varepsilon(t)$ .

As with Theorem 6, the performance bound on (KO) has three components: a bound on the generalization error, the bias due to optimizing with a scalar decision when the true decision is a function, and the finite-sample bias term, which is the same as in Theorems 5–6. The generalization error is of  $O(1/r_w(p)\sqrt{n})$ , so can be controlled by reducing  $r_w(p)$  by increasing the kernel bandwidth  $w$ . Setting  $w = O(\sqrt{p})$  gives an error which is of  $O(1/\sqrt{n})$ , which is as good as having the demand not depend on any features. When  $w$  is set to an arbitrarily large number,  $r_w(p) = 1$ , so the error rate  $O(1/\sqrt{n})$  cannot be improved upon. It is not surprising that the generalization error can be made small with large  $w$  since this corresponds to smoother comparisons of the feature vectors from the past to the one in period  $n+1$ . However, as with (NV-ERM2), increased  $w$  increases the second term, thus in practice  $w$  needs to be optimized over a reasonable range of values. Finally, the finite-sample bias term plays the same role as the corresponding terms in Theorems 5–6.

## 5. Case Study: Nurse Staffing in a Hospital Emergency Room

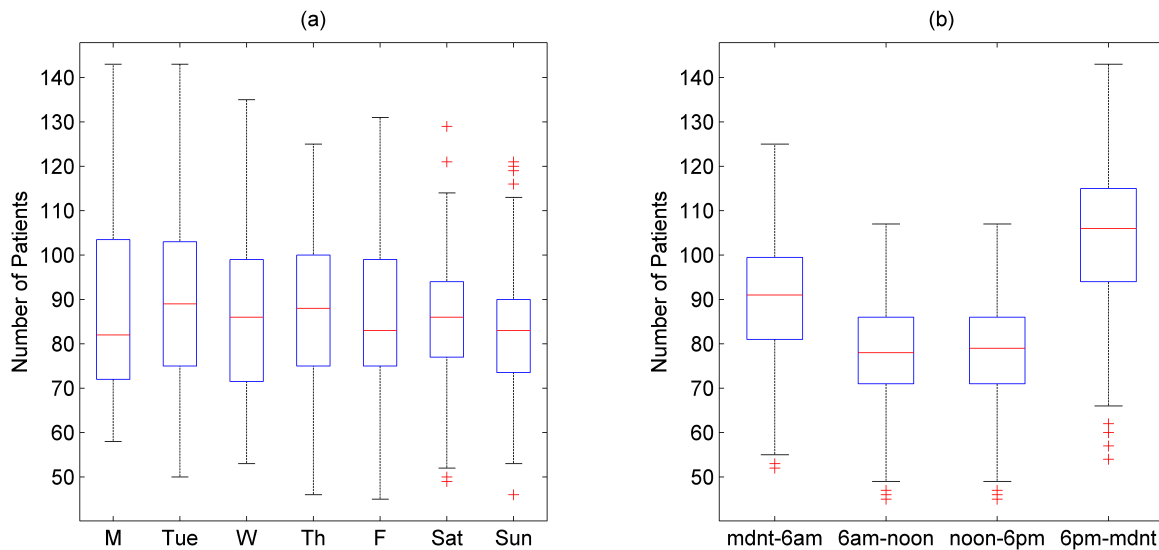
In this section, we compare the three algorithms introduced in Sec. 2, (NV-ERM1), (NV-ERM2) and (NV-KO) against the main data-driven benchmarks known in the literature and practice through an extensive empirical investigation. Although some analytical comparisons are possible under assumptions about the true demand model, the ultimate test of data-driven methods must be on real data sets. Furthermore, we report other practical observations such as the computational time required and trends in the optimal staffing solution.

In particular, we apply the three learning algorithms to find the optimal staffing levels of nurses for a hospital emergency room. As most hospitals in the developed world either impose or recommend a minimum nurse-to-patient ratio, we can approximate nurse staffing problem as a newsvendor problem if we assume the hospital incurs a linear underage cost if too many patients arrive and expensive agency nurses have to be called, and a linear overage cost if too many regular nurses are scheduled compared to the number of patients. As nurse staffing contributes to a significant portion of hospital operations (see e.g. Green et al. 2013) and as many hospitals are starting to harness the value of data, the nurse staffing problem is a natural setting to test the data-driven models introduced in this paper.



Our data comes from the emergency room of a large teaching hospital in the United Kingdom from July 2008 to June 2009. The data set includes the total number of patients in the emergency room at 2-hour intervals. We provide box plots of the number of patients by day and by time periods in Fig. 1. We assumed a nurse-to-patient ratio of 1 to 5, hence the demand is the total number of patients divided by 5. We do not require the staffing level to be an integer in our predictions, as multi-skilled workers could be used for part-time work. We also assumed that the hourly wage of an agency nurse is 2.5 times that of a regular nurse, that is  $b = 2.5/3.5$  and  $h = 1/3.5$ , resulting in a target fractile of  $r = b/(b + h) = 2.5/3.5$ . Although the exact agency nurse rate differs by location, experience and agency, our assumption is a modest estimate (Donnelly and Mulhern 2012).

We considered two sets of features: the first set being the day of the week, time of the day and  $m$  number of days of past demands; the second set being the first set plus the sample average of past demands and the differences in the order statistics of past demands, which is inspired by the analysis in Liyanage and Shanthikumar (2005) as described in Sec. 2.4.1. We refer to these features as *Operational Statistics* (OS) features. We used  $n = 1344$  past demand observations (16 weeks) as training data and computed the critical staffing level 3 periods ahead. We then recorded the out-of-sample newsvendor cost of the predicted staffing level on  $1344/2 = 672$  validation data on a rolling horizon basis, following the-rule-of thumb in Friedman et al. (2009) for choosing the size of the validation data set. Any parameter that needs calibration was calibrated on the validation data set. We then applied the algorithms to a test set of 672 unseen observations.



**Figure 1** A boxplot of the number of patients in the emergency room (a) by day and (b) by time period.

All computations were carried out on MATLAB2013a with the solver MOSEK and CVX, a package for specifying and solving convex programs (CVX Research 2012, Grant and Boyd 2008) on a Dell Precision T7600 workstation with two Intel Xeon E5-2643 processors, each of which has 4 cores, and 32.0 GB of RAM.

### 5.1. Methods considered

We investigate the following methods in detail:

1. SAA by day of the week: take a sample average of the training data set, by day of the week (since our training data set consists of 16 weeks of demand, there are  $1344/16 = 84$  observations for each day of the week). We note that this is reflective of nurse staffing done in practice.
2. Cluster + SAA: we take the vector of features, then first classify them into  $k = 2, \dots, 12$  clusters before applying SAA. This is an intuitive, and alternative method to use the feature data. For clustering, we use the k-means clustering algorithm.
3. Solve (NV-KO) with the Gaussian kernel, with the day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks) with and without Operational Statistics features. These constitute a total of two algorithms.
4. Solve (NV-ERM1) with the day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks) with and without Operational Statistics features, which were explained in Sec. 2.4.1. These constitute a total of two algorithms.
5. Solve (NV-ERM2) with the day of the week and time of the day features and 2 weeks of past demands with and without OS features for a range of regularization parameters. We investigate both  $\ell_1$  and  $\ell_2$  regularizations. These constitute a total of four algorithms.
6. Separated Estimation and Optimization (SEO): a common-sense approach to incorporating features in the newsvendor decision-making is by first regressing the demand on the features assuming a normally distributed error term (estimation) then applying the appropriate formula for the optimal quantile using the assumption of normality for the demand (optimization). We use day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks). We consider two cases, one without and one with OS features.
7. Separated Estimation and Optimization (SEO) approach with  $\ell_1$  or  $\ell_2$  regularization: we apply OLS regression with  $\ell_1$  or  $\ell_2$  or with no regularization to first estimate a demand model, then choosing the optimal quantile under the assumption of normally distributed demand. For the demand estimation step, we use day of the week and time of the day features and 2 weeks of past demands, and consider using and not using OS features.
8. We also consider Scarf's Minimax approach (Scarf et al. 1958).

In Table 2, we summarize the abbreviations used to describe the 16 different methods considered in this paper.

Abbreviation	Description	OS Features?	Reg.?	Free parameter
1a. SAA-day	SAA by day of the week	No	None	None
1b. Cluster+SAA	First cluster then SAA	No	no. of clusters	None
2a. Ker-0	solve (NV-KO) with Gaussian kernel	No	None	bandwidth
2b. Ker-OS	"	Yes	None	"
3a. NV-0	solve (NV-ERM1)	No	None	no. of days of past demand
3b. NV-OS	"	Yes	None	"
4a. NVreg1	solve (NV-ERM2)	No	Yes, $\ell_1$	regularization parameter
4b. NVreg1-OS	"	Yes	Yes, $\ell_1$	"
5a. NVreg2	"	No	Yes, $\ell_2$	"
5b. NVreg2-OS	"	Yes	Yes, $\ell_2$	"
6a. SEO-0	OLS regression + NV optimization	No	None	no. of days of past demand
6b. SEO-OS	"	No	None	"
7a. SEOreg1	Lasso regression + NV optimization	No	Yes, $\ell_1$	regularization parameter
7b. SEOreg1-OS	"	Yes	Yes, $\ell_1$	"
8a. SEOreg2	Ridge regression + NV optimization	No	Yes, $\ell_2$	"
8b. SEOreg2-OS	"	Yes	Yes, $\ell_2$	"
9. Scarf	Minimax optimization	No	None	no. of days of past demand

**Table 2** A summary of the methods considered.

## 5.2. Discussion of Results

In Table 3, we report the out-of-sample performance of the sixteen methods considered. We report the calibrated parameter (if any), the mean and the 95% confidence interval for the out-of-sample staffing cost in normalized units (parameters are calibrated by in-sample calculations, which can be found in Appendix E Tables EC.1–EC.7). In the last column, we report the annual cost savings of the method relative to SAA-day where there is a statistically significant net cost saving, assuming a regular nurse salary of £25,000 (which is the Band 4 nurse salary for the National Health Service in the United Kingdom in 2014) and standard working hours of 37.5 hours per week. Cost savings in USD are also reported, assuming an exchange rate of £1: USD 1.6.

The best result was obtained by the KO method with OS features, with bandwidth  $w = 1.62$ , which yields a cost improvement of 24% (a saving of £46,555 p.a.) relative to the best practice benchmark (“SAA-day”) with statistical significance at the 5% level. The next best results were obtained by the ERM method with  $\ell_1$  regularization and the KO method without OS features, which have average annual cost improvements of 23% (£44,219 p.a.) and 21% (£39,915 p.a.) respectively.

The computational costs of the feature-based methods are very different, however. The KO method is *three* orders of magnitude faster than the ERM-based methods. For instance, it takes just 0.0494 seconds to find the next optimal staffing level using the KO method with OS features, which is the best in terms of the out-of-sample cost, whereas it takes 114 seconds for the the ERM method

Method	Calibrated parameter	Avg. Comp. Time (per iteration)	Mean (95 % CI)	% savings rel. to SAA-day	Annual cost savings rel. to SAA-day
1a. SAA-day	—	14.0 s	1.523 ( $\pm 0.109$ )	—	—
1b. Cluster+SAA	—	14.9 s	1.424 ( $\pm 0.102$ )	—	—
2a. Ker-0	$w = 0.08$	0.0444 s	1.208 ( $\pm 0.146$ )	20.7%	£39,915 (\$ 63,864)
2b. Ker-OS	$w = 1.62$	0.0494 s	1.156 ( $\pm 0.140$ )	24.1%	£46,555 (\$ 74,488)
3a. NV-0	12 days	325 s	1.326 ( $\pm 0.100$ )	12.9%	£24,909 (\$ 39,854)
3b. NV-OS	4 days	360 s	1.463 ( $\pm 0.144$ )	—	—
4a. NVreg1	$1 \times 10^{-7}$	84.5 s	1.336 ( $\pm 0.100$ )	—	—
4b. NVreg1-OS	$1 \times 10^{-7}$	114 s	1.174 ( $\pm 0.113$ )	22.9%	£44,219 (\$ 70,750)
5a. NVreg2	$5 \times 10^{-7}$	79.6 s	1.336 ( $\pm 0.110$ )	—	—
5b. NVreg2-OS	$1 \times 10^{-7}$	107 s	1.215 ( $\pm 0.111$ )	20.2%	£39,065 (\$ 62,503)
6a. SEO-0	1 day	10.8 s	1.279 ( $\pm 0.099$ )	16.0%	£30,952 (\$ 49,523)
6b. SEO-OS	6 days	16.1 s	12.57 ( $\pm 10.63$ )	—	—
7a. SEOreg1	$5 \times 10^{-1}$	22.1 s	1.417 ( $\pm 0.106$ )	—	—
7b. SEOreg1-OS	$5 \times 10^{-3}$	25.9 s	11.95 ( $\pm 6.00$ )	—	—
8a. SEOreg2	$1 \times 10^{-1}$	26.6 s	1.392 ( $\pm 0.105$ )	—	—
8b. SEOreg2-OS	$5 \times 10^{-3}$	27.1 s	12.57 ( $\pm 10.63$ )	—	—
9. Scarf	12 days	20.8 s	1.593 ( $\pm 0.114$ )	—	—

**Table 3** A summary of results. We assume the hourly wage of an agency nurse is 2.5 times that of a regular nurse. We report the calibrated parameter (if any), the average computational time taken to solve one problem instance, and the mean and the 95% confidence interval for the out-of-sample staffing cost in normalized units. In the last column, we report the annual cost savings of the method relative to SAA-day where there is a statistically significant net cost saving, assuming a regular nurse salary of £25,000 (which is the Band 4 nurse salary for the National Health Service in the United Kingdom in 2014) and standard working hours. A dashed line represents cost differential that is not statistically significant. Cost savings in USD are also reported, assuming an exchange rate of £1: USD 1.6.

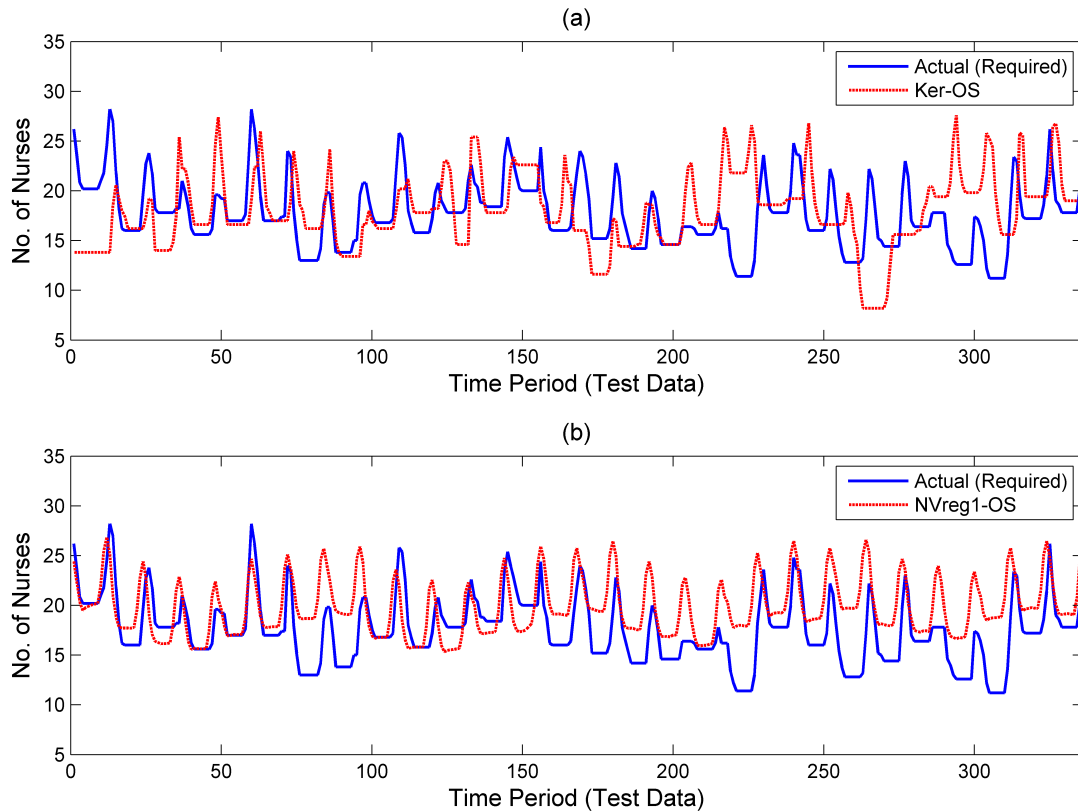
with  $\ell_1$  regularization, which is the second-best performing method. The KO method is also faster than SAA-day, SEO methods and Scarf by two orders of magnitude.

### 5.3. Optimal Staffing Decisions

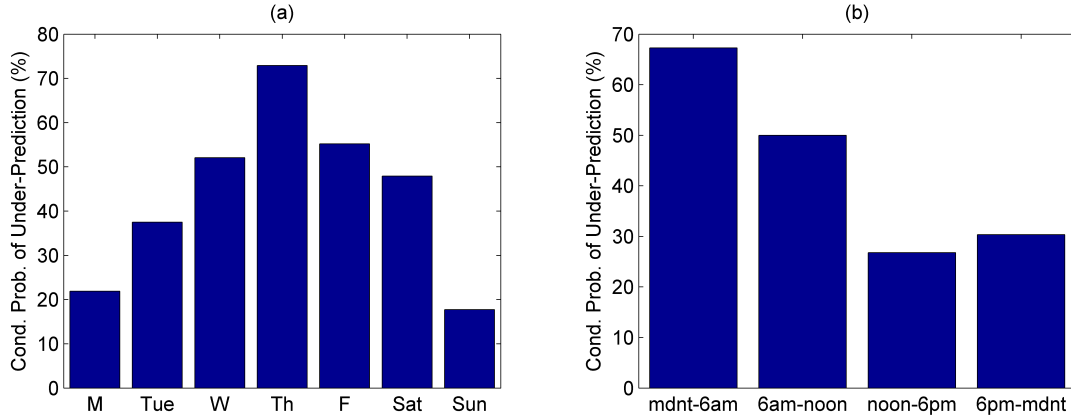
Let us further investigate the staffing decision of the best method: KO-OS with  $w = 1.62$ . In Fig. 2 (a), we display the staffing levels predicted by KO-OS with  $w = 1.62$  along with the actual required levels. For comparison, we also provide the staffing levels predicted by the second-best method, NVreg1-OS with  $\lambda = 1 \times 10^{-7}$  in Fig. 2 (b). A striking observation is that both KO-OS and NVreg1-OS methods anticipate periods of high demand fairly well, as evidenced by the matching of the peaks in the predicted and actual staffing levels. The two methods are otherwise quite different in the prediction; in particular, the KO-OS method balances both over-staffing and under-staffing, whereas NVreg1-OS method seems to systematically over-predict the staffing level.

Let us now suppose the hospital indeed implements our algorithm for its nurse staffing decisions. We wish to gain some insight into the predictions made by the algorithm. In particular, we would like to know when the hospital is over- or under-staffed, assuming the hospital chooses to implement the

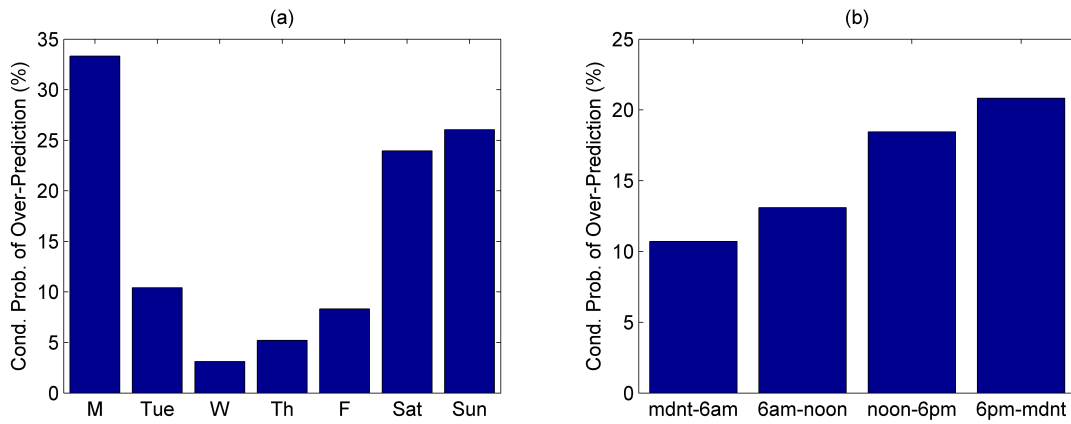
best possible method, provided by KO-OS with  $w = 1.62$ . In Figs. 3 and 4, we show the conditional probability (frequency) of under- and over-staffing by day of the week and by time period. We derive the following insights from these plots, which could be useful for patients and managers directly: (i) mid-week days are more likely to be under-staffed then weekends, thus, given the choice to visit the emergency room on a weekday or weekend, we would choose a weekend, (ii) the period from noon to midnight is substantially more likely to be over-staffed then the period from midnight to noon, thus, given the choice of time to visit the emergency room, we would choose visiting in the afternoon, and (iii) the algorithm is most likely to over-staff by at least 50% of the required level on a Monday then any other day of the week, hence, given the flexibility, we would choose to visit the emergency room on a Monday.



**Figure 2** (a) A time-series plot of actual staffing demand (solid blue) versus staffing levels predicted by KO-OS with  $w = 1.62$  (best method) on test data set in dotted red. (b) A time-series plot of actual staffing demand (solid blue) versus staffing levels predicted by NVreg1-OS with  $\lambda = 1 \times 10^{-7}$  (second best method) on test data set in dotted red.



**Figure 3** A plot of the conditional probabilities of under-staffing (a) by day and (b) by time period for KO-OS with  $w = 1.62$  (best method). The conditioning is done by the particular day or the time period, i.e. the probability of under-staffing given it is a Monday.



**Figure 4** A plot of the conditional probabilities of over-staffing by at least 50% (a) by day and (b) by time period for KO-OS with  $w = 1.62$  (best method). The conditioning is done by the particular day or the time period, i.e. the probability of over-staffing given it is a Monday.

## 6. Conclusion

We investigated the newsvendor problem when the DM has historical data on both demands and  $p$  features that are related to the demand. We have analyzed this problem using recent techniques from machine learning (algorithmic stability theory) as well as theoretical statistics (theory of quantile estimators). Rather than reiterate the contributions detailed in the Introduction, below we summarize the some practical insights that may not be obvious upon first reading (especially to readers unfamiliar with machine learning), and discuss potential directions for future research.

Some practical insights from this work are as follows. (i) There is not a single approach to solving the “Big Data” newsvendor problem. In this paper we proposed three approaches (ERM with and

without regularization, and KO) to using the feature-demand data set and compared with against a number of other potential approaches to solving the problem, with or without features. This is not to say these approaches are exhaustive — we are optimistic that new methods can be developed. (ii) A “Big Data”-driven decision-maker always needs to be weary of overfitting, where decisions can perform well in-sample but not so out-of-sample, hence leading to not only decisions that perform badly, but also decisions that are misleading. We have shown that overfitting can be controlled by a-priori feature selection or regularization, which gives the DM extra control to bias the decision in favour of improved out-of-sample generalization. (iii) What affects the true performance of the in-sample decision are the generalization error (a.k.a. overfitting error), bias from regularization and finite-sample bias simply from having a finite amount of data. The DM can control the generalization error and bias from regularization through the regularization parameter, and in practice needs to optimize over a range of parameter values on a separate validation data set. Finite-sample bias cannot be controlled except by collecting more data. (iv) While the KO method dominates all others in terms of the out-of-sample performance in our case study, this need not be the case for a different data set. We believe the most important point of our case study is in demonstrating how to carry out a careful data-driven investigation, not in the case-specific conclusion that the KO method performed the best. We note however the relative speed of the KO method is generalizable.

There are many directions for follow-up work, and we discuss a few. First, investigating how the dynamic inventory management problem can be solved with demand-feature data set remains an open question. Second, as mentioned in point (i) of the previous paragraph, the methods considered in this paper are not exhaustive — new methods can be developed, especially if different assumptions are made on the data-generating process than we have assumed here. It would be interesting to see, for instance, if Markov-modulated demand processes can be adapted to handle a large number of feature information. Finally, extending our theoretical results for more general classes of stochastic optimization problems remains an open task.

## Acknowledgments

This research was supported by London Business School Research and Material Development Scheme (Ban) and the National Science Foundation grant IIS-1053407 (Rudin). The authors thank Nicos Savva and Stefan Scholtes for providing the data for the case study. The authors are also grateful to Paul Zipkin, Chung-Piaw Teo, the anonymous referees and seminar attendees at a number of institutions and conferences for helpful suggestions.

## References

- Akaike, Hirotugu. 1974. A new look at the statistical model identification. *Automatic Control, IEEE Transactions on* **19**(6) 716–723.

- Arrow, Kenneth Joseph, Samuel Karlin, Herbert Scarf. 1958. *Studies in the mathematical theory of inventory and production*. 1, Stanford University Press.
- Azoury, Katy S. 1985. Bayes solution to dynamic inventory models under unknown demand distribution. *Management Science* **31**(9) 1150–1160.
- Belloni, Alexandre, Victor Chernozhukov. 2011.  $\ell_1$ -penalized quantile regression in high-dimensional sparse models. *The Annals of Statistics* **39**(1) 82–130.
- Bousquet, Olivier, André Elisseeff. 2002. Stability and generalization. *The Journal of Machine Learning Research* **2** 499–526.
- Burnetas, Apostolos N, Craig E Smith. 2000. Adaptive ordering and pricing for perishable products. *Operations Research* **48**(3) 436–443.
- Chaudhuri, Probal, et al. 1991. Nonparametric estimates of regression quantiles and their local bahadur representation. *The Annals of Statistics* **19**(2) 760–777.
- Chen, Xin, Melvyn Sim, Peng Sun. 2007. A robust optimization perspective on stochastic programming. *Operations Research* **55**(6) 1058–1071.
- Chernozhukov, Victor, Iván Fernández-Val, Alfred Galichon. 2010. Quantile and probability curves without crossing. *Econometrica* **78**(3) 1093–1125.
- Chernozhukov, Victor, Christian Hansen. 2008. Instrumental variable quantile regression: A robust inference approach. *Journal of Econometrics* **142**(1) 379–398.
- Csörgö, Miklos. 1983. *Quantile processes with statistical applications*. SIAM.
- CVX Research, Inc. 2012. CVX: Matlab software for disciplined convex programming, version 2.0. <http://cvxr.com/cvx>.
- Devroye, Luc, T Wagner. 1979a. Distribution-free inequalities for the deleted and holdout error estimates. *Information Theory, IEEE Transactions on* **25**(2) 202–207.
- Devroye, Luc, T Wagner. 1979b. Distribution-free performance bounds for potential function rules. *Information Theory, IEEE Transactions on* **25**(5) 601–604.
- Donnelly, L., M. Mulhern. 2012. Nhs pays £1,600 a day for nurses as agency use soars. <http://www.telegraph.co.uk/news/9400079/NHS-pays-1600-a-day-for-nurses-as-agency-use-soars.html>.
- Durrett, Rick. 2010. *Probability: theory and examples*. Cambridge university press.
- Feldman, Richard M. 1978. A continuous review (s, s) inventory system in a random environment. *Journal of Applied Probability* 654–659.
- Friedman, Jerome, Trevor Hastie, Robert Tibshirani. 2009. The elements of statistical learning: Data mining, inference, and prediction. *Springer Series in Statistics* .



- Gallego, Guillermo, Ilkyeong Moon. 1993. The distribution free newsboy problem: review and extensions. *Journal of the Operational Research Society* 825–834.
- Gallego, Guillermo, Özalp Özer. 2001. Integrating replenishment decisions with advance demand information. *Management Science* **47**(10) 1344–1360.
- Godfrey, Gregory A, Warren B Powell. 2001. An adaptive, distribution-free algorithm for the newsvendor problem with censored demands, with applications to inventory and distribution. *Management Science* **47**(8) 1101–1112.
- Grant, M., S. Boyd. 2008. Graph implementations for nonsmooth convex programs. V. Blondel, S. Boyd, H. Kimura, eds., *Recent Advances in Learning and Control*. Lecture Notes in Control and Information Sciences, Springer-Verlag Limited, 95–110. [http://stanford.edu/~boyd/graph\\_dcp.html](http://stanford.edu/~boyd/graph_dcp.html).
- Graves, Stephen C, Harlan C Meal, Sriram Dasu, Yuping Qui. 1986. Two-stage production planning in a dynamic environment. *Multi-stage production planning and inventory control*. Springer, 9–43.
- Green, Linda V, Sergei Savin, Nicos Savva. 2013. Nursevendor problem: Personnel staffing in the presence of endogenous absenteeism. *Management Science* **59**(10) 2237–2256.
- Hannah, Lauren, Warren Powell, David M Blei. 2010. Nonparametric density estimation for stochastic optimization with an observable state variable. *Advances in Neural Information Processing Systems*. 820–828.
- Heath, David C, Peter L Jackson. 1994. Modeling the evolution of demand forecasts ith application to safety stock analysis in production/distribution systems. *IIE transactions* **26**(3) 17–30.
- Hofmann, Thomas, Bernhard Schölkopf, Alexander J Smola. 2008. Kernel methods in machine learning. *The Annals of Statistics* 1171–1220.
- Huh, Woonghee Tim, Paat Rusmevichientong. 2009. A nonparametric asymptotic analysis of inventory planning with censored demand. *Mathematics of Operations Research* **34**(1) 103–123.
- Iida, Tetsuo, Paul H Zipkin. 2006. Approximate solutions of a dynamic forecast-inventory model. *Manufacturing & Service Operations Management* **8**(4) 407–425.
- Koenker, Roger. 2005. *Quantile regression*. Cambridge University Press.
- Kunnumkal, Sumit, Huseyin Topaloglu. 2008. Using stochastic approximation methods to compute optimal base-stock levels in inventory control problems. *Operations Research* **56**(3) 646–664.
- Levi, Retsef, Georgia Perakis, Joline Uichanco. 2015. The data-driven newsvendor problem: new bounds and insights. *Operations Research* **63**(6) 1294–1306.
- Levi, Retsef, Robin O Roundy, David B Shmoys. 2007. Provably near-optimal sampling-based policies for stochastic inventory control models. *Mathematics of Operations Research* **32**(4) 821–839.
- Liyanage, Liwan H, J George Shanthikumar. 2005. A practical inventory control policy using operational statistics. *Operations Research Letters* **33**(4) 341–348.

- Lovejoy, William S. 1990. Myopic policies for some inventory models with uncertain demand distributions. *Management Science* **36**(6) 724–738.
- Lovejoy, William S. 1992. Stopped myopic policies in some inventory models with generalized demand processes. *Management Science* **38**(5) 688–707.
- Lu, Xiangwen, Jing-Sheng Song, Amelia Regan. 2006. Inventory planning with forecast updates: Approximate solutions and cost error bounds. *Operations Research* **54**(6) 1079–1097.
- Manton, Jonathan H, Pierre-Olivier Amblard, et al. 2015. A primer on reproducing kernel hilbert spaces. *Foundations and Trends® in Signal Processing* **8**(1-2) 1–126.
- Nadaraya, Elizbar A. 1964. On estimating regression. *Theory of Probability & Its Applications* **9**(1) 141–142.
- Parzen, Emanuel. 1979. Nonparametric statistical data modeling. *Journal of the American Statistical Association* **74**(365) 105–121.
- Perakis, Georgia, Guillaume Roels. 2008. Regret in the newsvendor model with partial information. *Operations Research* **56**(1) 188–203.
- Powell, Warren, Andrzej Ruszczyński, Huseyin Topaloglu. 2004. Learning algorithms for separable approximations of discrete stochastic optimization problems. *Mathematics of Operations Research* **29**(4) 814–836.
- Rockafellar, R Tyrell. 1997. *Convex analysis*, vol. 28. Princeton University Press.
- Rogers, William H, Terry J Wagner. 1978. A finite sample distribution-free performance bound for local discrimination rules. *The Annals of Statistics* 506–514.
- Scarf, Herbert. 1959a. Bayes solutions of the statistical inventory problem. *The Annals of Mathematical Statistics* 490–508.
- Scarf, Herbert. 1959b. *The optimality of (s,S) policies in the dynamic inventory problem. Mathematical Methods in the Social Science, KJ Arrow, S. Karlin, P. Suppes, eds..* Stanford University Press, Stanford.
- Scarf, Herbert, KJ Arrow, S Karlin. 1958. A min-max solution of an inventory problem. *Studies in the Mathematical Theory of Inventory and Production* **10** 201–209.
- Schwarz, Gideon. 1978. Estimating the dimension of a model. *The Annals of Statistics* **6**(2) 461–464.
- See, Chuen-Teck, Melvyn Sim. 2010. Robust approximation to multiperiod inventory management. *Operations Research* **58**(3) 583–594.
- Shapiro, Alexander, Darinka Dentcheva, Andrzej P Ruszczyński. 2009. *Lectures on stochastic programming: modeling and theory*, vol. 9. SIAM.
- Song, Jing-Sheng, Paul Zipkin. 1993. Inventory control in a fluctuating demand environment. *Operations Research* **41**(2) 351–370.

- Takeuchi, Ichiro, Quoc V Le, Timothy D Sears, Alexander J Smola. 2006. Nonparametric quantile estimation. *The Journal of Machine Learning Research* **7** 1231–1264.
- Vapnik, Vladimir N. 1998. *Statistical learning theory*. Wiley.
- Watson, Geoffrey S. 1964. Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A* 359–372.

# Electronic Companion to “The Big Data Newsvendor: Practical Insights from Machine Learning”

## Appendix A: Solving (NV-ERM) over a Nonlinear Function Class

We can solve (NV-ERM) over a nonlinear function class by the use of kernels. Kernels allow to mapping of the problem (NV-ERM) to a higher dimensional space (known as a reproducing kernel Hilbert space) in a convenient way. The computations are approximately as difficult as in the original feature space, but the order quantity can now be a highly nonlinear function of the original features.

Consider the feature vector  $\mathbf{x}$  from the original feature space  $\mathcal{X} \subset \mathbb{R}^p$ . Suppose we transform this to a new feature,  $\phi(\mathbf{x})$ , of possibly infinite dimensions.  $\phi(\mathbf{x})$  may be a vector of dimension greater than  $p$ , or a function. By assuming  $\phi(\mathbf{x})$  belongs in a reproducing kernel Hilbert space with a known kernel  $k(\cdot, \cdot)$ , we can compute inner products  $\langle \phi(\mathbf{x}_1), \phi(\mathbf{x}_2) \rangle = k(\mathbf{x}_1, \mathbf{x}_2)$  by the reproducing property. As long as we can compute kernels between any two elements  $\mathbf{x}_1, \mathbf{x}_2$ , we can operate in this higher dimensional space, without having to specify the map  $\phi(\cdot)$ .

We can thus transform the newsvendor problem (NV-ERM) to optimize over a class of nonlinear decisions by transforming  $\mathbf{x}$  to  $\phi(\mathbf{x})$  via a kernel. Let  $q$  be a function within the reproducing kernel Hilbert space and  $q_0$  a scalar decision variable. Then the primal optimization problem becomes:

$$\begin{aligned} \min_{q, u_i, o_i} \quad & \frac{1}{n} \sum_{i=1}^n (bu_i + ho_i) + \lambda \|q\|_2^2 \\ \text{s.t. } \forall i = 1, \dots, n: \quad & u_i \geq d_i - q_0 - \langle \phi(\mathbf{x}_i), q \rangle \\ & o_i \geq q_0 + \langle \phi(\mathbf{x}_i), q \rangle - d_i \\ & u_i, o_i \geq 0, \end{aligned} \tag{NV-ERM2-KER}$$

where the optimal order quantity will be  $q_{\text{opt}}(\mathbf{x}) = \langle \phi(\mathbf{x}), q^* \rangle + q_0^*$ , where  $q^*$  and  $q_0^*$  is the solution to the optimization problem. Following the standard approaches we can derive the convex dual, using notation  $\mathbf{K}$  as the matrix of  $k(\mathbf{x}_i, \mathbf{x}_l)$  values, and using  $\boldsymbol{\alpha}$  as the vector of  $n$  combined Lagrange multipliers. Note that if  $k$  is a valid kernel (i.e. an inner product for a reproducing kernel Hilbert space), the gram matrix  $\mathbf{K}$  is positive semi-definite. The dual problem is thus:

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} - \boldsymbol{\alpha}^T \mathbf{d} \\ \text{s.t. } \forall i = 1, \dots, n: \quad & \frac{-h}{2\lambda n} \leq \alpha_i \leq \frac{b}{2\lambda n}, \text{ and} \\ & \sum_i \alpha_i = 0. \end{aligned}$$

If  $\alpha^*$  is the solution to this dual problem, which we can compute with a quadratic programming solver, the solution to the primal will be given by:

$$q^*(\mathbf{x}) = \sum_i \alpha_i^* \phi(\mathbf{x}_i),$$

and  $q_0^*$  is  $d_i - \langle \phi(\mathbf{x}_i), q^*(\mathbf{x}_i) \rangle = d_i - \sum_{l=1}^n \alpha_l^* k(\mathbf{x}_i, \mathbf{x}_l)$  for any  $i$  where  $\alpha_i^*$  is neither  $\frac{-h}{2\lambda n}$  nor  $\frac{b}{2\lambda n}$ . This means that the optimal order quantity  $q(\mathbf{x})$  for any new  $\mathbf{x}$  can be evaluated directly, using only computations of  $k(\mathbf{x}, \mathbf{x}_i)$ . This is given by:

$$q_{\text{opt}}(\mathbf{x}) = \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + q_0^*.$$

One can choose any kernel  $k$  so long as it satisfies the property of being a reproducing kernel, without calculating (or even knowing) the map  $\phi$  to the reproducing kernel Hilbert space. Typical kernels are polynomial kernels or Gaussian kernels. The kernel must be symmetric, as it needs to be an inner product in the reproducing kernel Hilbert space, and inner products are symmetric. Even if the kernel function is much more complicated than a standard inner product,  $q_{\text{opt}}(\mathbf{x})$  can be evaluated for any  $\mathbf{x}$ , by solving the dual formulation and using the formulae above. This is the generalization of (NV-ERM) to nonlinear decision rules.

For further readings on the topic, we direct the readers to [Hofmann et al. \(2008\)](#) and [Manton et al. \(2015\)](#).

## Appendix B: Proofs of results in Sec. 3

*Proof of Lemma 1:* The feature-based algorithm (NV-ERM1) solves

$$\begin{aligned} \min_{q(x)=q^0(1-x)+q^1x} \hat{R}(q(x); S_n) &= \frac{1}{n} \sum_{i=1}^n [b(d_i(x) - q(x))^+ + h(q(x) - d_i(x))^+] \\ &= \min_{q(x)=q^0(1-x)+q^1x} \frac{1}{n_0} \sum_{i:x_i=0} [b(d_i^0 - q^0)^+ + h(q^0 - d_i^0)^+] + \frac{1}{n_1} \sum_{i:x_i=1} [b(d_i^1 - q^1)^+ + h(q^1 - d_i^1)^+] \\ &= \min_{q^0 \geq 0} \left\{ \frac{1}{n_0} \sum_{i:x_i=0} [b(d_i^0 - q^0)^+ + h(q^0 - d_i^0)^+] \right. \\ &\quad \left. + \min_{q^1 \geq 0} \left\{ \frac{1}{n_1} \sum_{i:x_i=1} [b(d_i^1 - q^1)^+ + h(q^1 - d_i^1)^+] \right\} \right\}, \end{aligned} \tag{EC.1}$$

where the outer and inner minimization problems correspond to the SAA problem for the subsample of data corresponding to  $x = 0$  and  $x = 1$  respectively. Hence the solutions are the corresponding SAA solutions for the appropriate subsample of data, which is the well-known critical fractile of the inverse empirical cdf as in (4).  $\square$

*Proof of Theorem 1:* Under Condition A, the following strong result holds via Theorem 4.1.2. pp. 31 of Csörgö (1983): there exists, for each  $n_i$ ,  $i = 0, 1$ , a Brownian Bridge  $\{B_{n_i}(y), 0 \leq y \leq 1\}$  such that

$$\sup_{0 < y < 1} \left| f_i(F_i^{-1}(y))(\hat{F}_i^{-1}(y) - F_i^{-1}(y)) - \frac{B_{n_i}(y)}{\sqrt{n_i}} \right| \stackrel{a.s.}{=} O\left(\frac{\log n_i}{n_i}\right). \quad (\text{EC.2})$$

The above implies, for  $y = r$ :

$$\begin{aligned} & \left| (\hat{F}_i^{-1}(r) - F_i^{-1}(r)) - \frac{B_{n_i}(r)}{f_i(F_i^{-1}(r))\sqrt{n_i}} \right| \stackrel{a.s.}{\leq} O\left(\frac{\log n_i}{n_i}\right) \\ \implies & \left| \hat{F}_i^{-1}(r) - F_i^{-1}(r) \right| \stackrel{a.s.}{\leq} \frac{B_{n_i}(r)}{f_i(F_i^{-1}(r))\sqrt{n_i}} + O\left(\frac{\log n_i}{n_i}\right) \\ \implies & \left| \mathbb{E}[\hat{F}_i^{-1}(r)] - F_i^{-1}(r) \right| \leq \mathbb{E} \left| \hat{F}_i^{-1}(r) - F_i^{-1}(r) \right| \leq \frac{\mathbb{E}B_{n_i}(r)}{f_i(F_i^{-1}(r))\sqrt{n_i}} + O\left(\frac{\log n_i}{n_i}\right) = O\left(\frac{\log n_i}{n_i}\right), \end{aligned}$$

where the last line uses Jensen's inequality and the fact that the mean of a Brownian Bridge is zero everywhere. Hence we get both the finite-sample bias result and the asymptotic optimality result.  $\square$

*Proof of Lemma 2:* This is equivalent to the SAA solution for the complete data set.  $\square$

*Proof of Theorem 2:* Proof of (15) parallels that of Proposition 1. Proof of (16) then follows from (15) and Proposition 1.

Proof of (17): The asymptotic convergence of  $\hat{q}_n^{SAA}$  to its true value is again due to the asymptotic convergence of the sample quantile estimator, as shown in Theorem 1. The statement then follows.  $\square$

*Proof of Theorem 3:* The SAA decision is given by

$$\begin{aligned} \hat{q}_n^{SAA}(\tilde{\mathbf{x}}) &= (\hat{F}_n^0)^{-1} \left( \frac{b}{b+h} \right) + \bar{d}_n \\ &= (\hat{F}_n^0)^{-1} \left( \frac{b}{b+h} \right) + \frac{1}{n} \sum_{i=1}^n \beta^\top \mathbf{X}_i + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \\ &\stackrel{a.s.}{\rightarrow} Q_\varepsilon \left( \frac{b}{b+h} \right) + \beta^\top \mathbb{E}[\mathbf{X}] + 0, \end{aligned}$$

where the convergence of the first term is due to the Glivenko-Cantelli theorem and the convergence of the last two terms is due to the Strong Law of Large Numbers (SLLN) (see Durrett 2010).  $\square$

*Proof of Theorem 4:* Let  $\underline{D} \leq q_1 < q_2 \leq \bar{D}$  be two newsvendor decisions. Then by considering the difference  $C(q_2; \cdot) - C(q_1; \cdot)$  for three different regions:  $D < q_1$ ,  $q_1 \leq D \leq q_2$  and  $D > q_2$ , one can show

$$\begin{aligned} C(q_2; D) - C(q_1; D) &= -b(q_2 - q_1)\mathbb{I}(D > q_2) + h(q_2 - q_1)\mathbb{I}(D < q_1) + [bq_1 + hq_2 - (b+h)D]\mathbb{I}(q_1 \leq D \leq q_2) \\ &= -b(q_2 - q_1)\mathbb{I}(D > q_2) + h(q_2 - q_1)\mathbb{I}(D < q_1) + [b(q_2 - (q_2 - q_1)) + hq_2 - (b+h)D]\mathbb{I}(q_1 \leq D \leq q_2) \\ &= -b(q_2 - q_1)\mathbb{I}(D > q_1) + h(q_2 - q_1)\mathbb{I}(D < q_1) + (b+h)(q_2 - D)\mathbb{I}(q_1 \leq D \leq q_2) \end{aligned} \quad (\text{EC.3})$$

Also,

$$\begin{aligned}
C(q_2; D) - C(q_1; D) &= -b(q_2 - q_1)\mathbb{I}(D > q_2) + h(q_2 - q_1)\mathbb{I}(D < q_1) + [bq_1 + hq_2 - (b + h)D]\mathbb{I}(q_1 \leq D \leq q_2) \\
&= -b(q_2 - q_1)\mathbb{I}(D > q_2) + h(q_2 - q_1)\mathbb{I}(D < q_1) + [bq_1 + h(q_1 + (q_2 - q_1)) - (b + h)D]\mathbb{I}(q_1 \leq D \leq q_2) \\
&= -b(q_2 - q_1)\mathbb{I}(D > q_2) + h(q_2 - q_1)\mathbb{I}(D < q_2) - (b + h)(D - q_1)\mathbb{I}(q_1 \leq D \leq q_2) \tag{EC.4}
\end{aligned}$$

Case I: Let us suppose  $q_1 = q^*$ . Then  $\mathbb{P}(D < q_1) = b/(b + h)$  and  $\mathbb{P}(D > q_1) = h/(b + h)$ , since we assume the demand is continuous. Thus the expectation of the first two terms in (EC.3) is equal to zero. The expected cost difference then becomes

$$|\mathbb{E}C(q_2; D) - \mathbb{E}C(q_1; D)| = (b + h)\mathbb{E}[(q_2 - D)\mathbb{I}(q_1 \leq D \leq q_2)],$$

which clearly increases as  $q_2$  increases (i.e. deviates further from  $q_1$ ).

Case II: Let us suppose  $q_2 = q^*$ . Then  $\mathbb{P}(D < q_2) = b/(b + h)$  and  $\mathbb{P}(D > q_2) = h/(b + h)$ , since we assume the demand is continuous. Thus the expectation of the first two terms in (EC.4) is equal to zero. The expected cost difference then becomes

$$|\mathbb{E}C(q_2; D) - \mathbb{E}C(q_1; D)| = (b + h)\mathbb{E}[(D - q_1)\mathbb{I}(q_1 \leq D \leq q_2)],$$

which clearly increases as  $q_1$  decreases (i.e. deviates further from  $q_2$ ).

In sum, the expected cost difference between  $q^*$  and  $\hat{q}$  increases as  $|q^* - \hat{q}|$  increases.  $\square$

## Appendix C: Proofs of Theorems in Sec. 4

We start by first showing that the optimal (in-sample) decisions obtained by solving the three optimization problems are *algorithmically stable*, which means that the true expected cost of the decision is not sensitive to changes in the data set. The stability of a decision is thus a desirable property and is intimately linked to how well the decision performs on new observations. We then quantify the performance of the three decisions on a new observation in the form of *generalization bounds*, and extend them to quantify how well the in-sample decisions perform relative to the true optimal decision in terms of the expected out-of-sample cost.

Generalization bounds are probabilistic bounds on the out-of-sample performance of in-sample predictions, and are useful in highlighting quantities that are important in prediction. There are several types of generalization bounds that are known in the statistical learning theory literature. The type of bounds we employ are *stability* bounds, because as algorithm-specific bounds, they generate the most insightful results for the feature-based newsvendor problem. In contrast, *uniform* generalization bounds, which are more common in statistical learning theory and perhaps more

familiar to the reader, are algorithm-independent. Uniform generalization bounds do not consider the way in which the algorithm searches the space of possible models, and as such they lack specific insights about prediction.

Stability bounds can be thought of as a form of Hoeffding’s inequality for algorithms (recall, for instance, that the sample-size bounds of [Levi et al. 2007](#) and [Levi et al. 2015](#) are critically based on Hoeffding’s inequality). To apply Hoeffding’s inequality, we need to know a bound on the values of a random variable, whereas to apply stability bounds, we need to know a bound on the stability of an algorithm to a random data set. The challenging part is to prove a stability result that is as strong as possible for the algorithm, which we here achieve for the newsvendor problem. Stability bounds have origins in the 1970’s ([Rogers and Wagner 1978](#), [Devroye and Wagner 1979a,b](#)), and recent work includes that of [Bousquet and Elisseeff \(2002\)](#).

To show that our newsvendor algorithms are stable, we first show that the newsvendor cost on the training set does not change very much when one of the training examples changes. Stability of a decision is a desirable property, and these bounds show that stability is intimately linked to how well the decision performs on new observations. One of the main contributions of this section is thus in showing that (NV-ERM1), (NV-ERM2) and (NV-KO) are strongly stable, which are important results in their own right.

Our algorithms for the learning newsvendor problem turn out to have a very strong stability property, namely it is *uniformly stable*. We define stability and uniform stability below.

DEFINITION EC.1 (UNIFORM STABILITY, [BOUSQUET AND ELISSEEFF \(2002\)](#) DEF 6 PP. 504).

A symmetric algorithm  $A$  has uniform stability  $\alpha$  with respect to a loss function  $\ell$  if for all  $S_n \in \mathcal{Z}^n$  and for all  $i \in \{1, \dots, n\}$ ,

$$\|\ell(A_{S_n}, \cdot) - \ell(A_{S_n^i}, \cdot)\|_\infty \leq \alpha. \quad (\text{EC.5})$$

Furthermore, an algorithm is *uniformly stable* if  $\alpha = \alpha_n \leq O(1/n)$ .

In what follows, the random demand is denoted by  $D$ , and is assumed to be bounded:  $D \in \mathcal{D} := [0, \bar{D}]$ .

As before, the historical (‘training’) set of data is given by  $S_n = \{(\mathbf{x}_i, d_i)\}_{i=1}^n$ .

We start with the results for (NV-ERM1).

PROPOSITION EC.1 (**Uniform stability of (NV-ERM1)**). *The learning algorithm (NV-ERM1) with iid data is symmetric and uniformly stable with respect to the newsvendor cost function  $C(\cdot, \cdot)$  with stability parameter*

$$\alpha_n = \frac{\bar{D}(b \vee h)^2}{(b \wedge h)} \frac{p}{n}. \quad (\text{EC.6})$$



Here the notation  $b \vee h$  indicates the maximum value of  $b$  and  $h$ , and  $b \wedge h$  indicates the minimum of the two. This bound illuminates that if one of  $b$  or  $h$  is too large and the other too small, in other words if the backordering and holding costs are highly asymmetric, then the algorithm is very sensitive to small changes in the data set. The algorithm is more stable when the target order quantity is closer to the median of the demand distribution.

The stability result, coupled with a Hoeffding/McDiarmid based argument (in our case, we are using more contemporary versions of the results due to [Bousquet and Elisseeff 2002](#)), yields the following:

**PROPOSITION EC.2 (Generalization Bound for (NV-ERM1)).** *Let  $\hat{q}$  be the model produced by Algorithm (NV-ERM1). The following bound holds with probability at least  $1 - \delta$  over the random draw of the sample  $S_n$ , where each element of  $S_n$  is drawn iid from an unknown distribution on  $\mathcal{X} \times \mathcal{D}$ :*

$$\frac{|R_{true}(\hat{q}) - \hat{R}(\hat{q}; S_n)|}{(b \vee h)\bar{D}} \leq \frac{2(b \vee h)}{b \wedge h} \frac{p}{n} + \left( \frac{4(b \vee h)}{b \wedge h} p + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}}.$$

For large  $p/n$ , we had suggested finding the optimal order quantity by solving (NV-ERM2) instead. In this case, the generalization is driven by the regularization parameter, rather than the ratio  $p/n$ , as well as  $b$  and  $h$  as before. First is the stability result.

**PROPOSITION EC.3 (Uniform stability of (NV-ERM2)).** *The learning algorithm (NV-ERM2) is symmetric, and is uniformly stable with respect to the newsvendor cost function  $C$  with stability parameter*

$$\alpha_n^r = \frac{(b \vee h)^2 X_{\max}^2 p}{2n\lambda}. \quad (\text{EC.7})$$

The new stability parameter is  $O(p/n\lambda)$ . Thus the stability of the algorithm can be controlled via the regularization parameter  $\lambda$ .

We will use the following lemma in the proof of Propositions [EC.1](#) and [EC.5](#).

**LEMMA EC.1 (Tight uniform bound on the newsvendor cost).** *The newsvendor cost function  $C(\cdot, \cdot)$  is bounded by  $(b \vee h)\bar{D}$ , which is tight in the sense that:*

$$\sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} |C(q, D(\mathbf{x}))| = \bar{D}(b \vee h).$$

*Proof of Lemma EC.1:* Clearly,  $\bar{D}(b \vee h)$  is an upper bound on  $|C(q, d)|$  for all  $q, d \in [0, \bar{D}]$ . Now if  $d = 0$  and  $q = \bar{D}$ ,  $|C(q, d)| = \bar{D}h$ . Conversely, if  $d = \bar{D}$  and  $q = 0$ ,  $|C(q, d)| = \bar{D}b$ . Hence the upper bound is attained.  $\square$

Now for the proof of the Proposition [EC.1](#).

*Proof of Proposition EC.1:* Symmetry follows from the fact that the data-generating process is iid. For stability, we will change our notation slightly to make the dependence on  $n$  and  $S_n$  explicit. Let

$$q_n(\mathbf{x}) := \mathbf{q}_n^\top \mathbf{x} = \sum_{j=1}^p q_n^j x_j$$

and

$$q_{n \setminus i}(\mathbf{x}) := \mathbf{q}_{n \setminus i}^\top \mathbf{x} = \sum_{j=1}^p q_{n \setminus i}^j x_j$$

where

$$[q_n^1, \dots, q_n^p] = \arg \min_{\mathbf{q}=[q^1, \dots, q^p]} \hat{R}(\mathbf{q}; S_n) = \frac{1}{n} \sum_{j=1}^n \left[ b \left( d_j - \sum_{j=1}^p q^j x_j \right)^+ + h \left( \sum_{j=1}^p q^j x_j - d_j \right)^+ \right]$$

is the solution to (NV-ERM1) for the set  $S_n$ , and

$$(q_{n \setminus i}^1, q_{n \setminus i}^p) = \arg \min_{\mathbf{q}=[q^1, \dots, q^p]} \hat{R}(\mathbf{q}; S_n^{\setminus i}) = \frac{1}{n} \sum_{j=1}^n \left[ b \left( d_j - \sum_{j=1}^p q^j x_j \right)^+ + h \left( \sum_{j=1}^p q^j x_j - d_j \right)^+ \right]$$

is the solution to (NV-ERM1) for the set  $S_n^{\setminus i}$ . Note that:

$$\hat{R}(\mathbf{q}; S_n) = \frac{n-1}{n} \hat{R}(\mathbf{q}; S_n^{\setminus i}) + \frac{1}{n} \tilde{R}(\mathbf{q}; S_i),$$

where  $S_i = (\mathbf{x}_i, d_i)$ .

For a fixed  $\mathbf{x}$ , we have, by the Lipschitz property of  $C(q; \cdot)$ ,

$$\sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} |C(q_n(\mathbf{x}), D(\mathbf{x})) - C(q_{n \setminus i}(\mathbf{x}), D(\mathbf{x}))| \leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} (b \vee h) |q_n(\mathbf{x}) - q_{n \setminus i}(\mathbf{x})|.$$

So we want to bound

$$|q_n(\mathbf{x}) - q_{n \setminus i}(\mathbf{x})| = \left| \sum_{j=1}^p q_n^j x_j - \sum_{j=1}^p q_{n \setminus i}^j x_j \right|.$$

By the convexity of the function  $\hat{R}_n(\cdot, S)$ , we have (see Section 23 of Rockafellar (1997)):

$$\sum_{j=1}^p \nu_j (q_{n \setminus i}^j - q_n^j) \leq \hat{R}(\mathbf{q}_{n \setminus i}; S_n) - \hat{R}(\mathbf{q}_n; S_n)$$

for all  $\boldsymbol{\nu} = [\nu_1, \dots, \nu_m] \in \partial \hat{R}(q_n; S_n)$  (set of subgradients of  $\hat{R}(\cdot, S_n)$  at  $q_n$ ). Furthermore, because  $0 \in \partial \hat{R}(q_n; S_n)$  by the optimality of  $q_n$ , we have

$$0 \leq \max_{\boldsymbol{\nu} \in \partial \hat{R}(q_n; S_n)} \sum_{j=1}^p \nu_j (q_{n \setminus i}^j - q_n^j) \leq \hat{R}(\mathbf{q}_{n \setminus i}; S_n) - \hat{R}(\mathbf{q}_n; S_n)$$

where the max over  $\boldsymbol{\nu}$  can be attained because  $\partial\hat{R}(q_n; S_n)$  is a compact set. Denote this maximum  $\boldsymbol{\nu}^*$ . We thus have

$$\begin{aligned}\hat{R}(\mathbf{q}_{n\setminus i}; S_n) - \hat{R}(\mathbf{q}_n; S_n) &\geq |\boldsymbol{\nu}^{*\top}(\mathbf{q}_{n\setminus i} - \mathbf{q}_n)| = \sum_{j=1}^p \nu_j^* (q_{n\setminus i}^j - q_n^j) \\ &\geq |\nu_j^* (q_{n\setminus i}^j - q_n^j)| = |\nu_j^*| |q_{n\setminus i}^j - q_n^j| \quad \text{for all } j = 1, \dots, p\end{aligned}$$

where the second inequality is because  $\nu_j^* (q_{n\setminus i}^j - q_n^j) > 0$  for all  $j$  because  $\hat{R}(\cdot; S_n)$  is piecewise linear and nowhere flat. Thus we get, for all  $j = 1, \dots, p$ ,

$$|q_{n\setminus i}^j - q_n^j| \leq \frac{\hat{R}(\mathbf{q}_{n\setminus i}; S_n) - \hat{R}(\mathbf{q}_n; S_n)}{|\nu_j^*|}.$$

Let us bound  $\hat{R}(\mathbf{q}_{n\setminus i}; S_n) - \hat{R}(\mathbf{q}_n; S_n)$ . Note

$$\begin{aligned}\hat{R}(\mathbf{q}_n; S_n) &= \frac{n-1}{n} \hat{R}(\mathbf{q}_n; S_n^{\setminus i}) + \frac{1}{n} \hat{R}(\mathbf{q}_n; S_i) \\ &\geq \frac{n-1}{n} \hat{R}(\mathbf{q}_{n\setminus i}; S_n^{\setminus i})\end{aligned}$$

since  $\mathbf{q}_{n\setminus i}$  is the minimizer of  $\hat{R}(\cdot; S_n^{\setminus i})$ . Also,  $\hat{R}(\mathbf{q}_n; S_n) \leq \hat{R}(\mathbf{q}_{n\setminus i}; S_n)$  since  $q_n$  is by definition the minimizer of  $\hat{R}(\cdot; S_n)$ . Putting these together, we get

$$\begin{aligned}\frac{n-1}{n} \hat{R}(\mathbf{q}_{n\setminus i}; S_n^{\setminus i}) - \hat{R}(\mathbf{q}_{n\setminus i}; S_n) &\leq \hat{R}(\mathbf{q}_n; S_n) - \hat{R}(\mathbf{q}_{n\setminus i}; S_n) \leq 0 \\ \implies |\hat{R}(\mathbf{q}_n; S_n) - \hat{R}(\mathbf{q}_{n\setminus i}; S_n)| &\leq \left| \frac{n-1}{n} \hat{R}(\mathbf{q}_{n\setminus i}; S_n^{\setminus i}) - \hat{R}(\mathbf{q}_{n\setminus i}; S_n) \right| \\ &= \left| \frac{n-1}{n} \hat{R}(\mathbf{q}_{n\setminus i}; S_n^{\setminus i}) - \frac{n-1}{n} \hat{R}(\mathbf{q}_{n\setminus i}; S_n^{\setminus i}) - \frac{1}{n} \hat{R}(\mathbf{q}_{n\setminus i}; S_i) \right| \\ &= \frac{1}{n} |\hat{R}(\mathbf{q}_{n\setminus i}; S_i)|.\end{aligned}$$

Thus

$$\begin{aligned}\sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} (b \vee h) |q_n(\mathbf{x}) - q_{n\setminus i}(\mathbf{x})| &\leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} (b \vee h) \left( \sum_{j=1}^p |q_n^j - q_{n\setminus i}^j| |x_j| \right) \\ &\leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} (b \vee h) \cdot \sum_{j=1}^p \frac{|x_j|}{|\nu_j^*|} \cdot (\hat{R}(\mathbf{q}_{n\setminus i}; S_n) - \hat{R}(\mathbf{q}_n; S_n)) \\ &\leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} \frac{b \vee h}{n} \cdot \sum_{j=1}^p \frac{|x_j|}{|\nu_j^*|} \cdot |\hat{R}(\mathbf{q}_{n\setminus i}; S_i)|.\end{aligned}\tag{EC.9}$$

We can further simplify the upper bound (EC.9) as follows. Recall that  $\boldsymbol{\nu}^*$  is the subgradient of  $\hat{R}(\cdot; S_n)$  at  $\mathbf{q}_n$  that maximizes  $\sum_{j=1}^p \nu_j (q_{n\setminus i}^j - q_n^j)$ ; and as  $\partial\hat{R}(\mathbf{q}_n; S_n)$  is compact (by the convexity of  $\hat{R}(\cdot; S_n)$ ), we can compute  $\boldsymbol{\nu}^*$  exactly. It is straightforward to show:

$$\nu_j^* = \begin{cases} -bx_j & \text{if } q_{n\setminus i}^j - q_n^j \leq 0 \\ hx_j & \text{if } q_{n\setminus i}^j - q_n^j \geq 0 \quad \forall j. \end{cases}$$

We can thus bound  $1/|\nu_j^*|$  by  $1/[(b \wedge h)|x_j|]$ . By using the tight uniform upper bound  $(b \vee h)\bar{D}$  on each term of  $|\hat{R}(\cdot, \cdot)|$  from Lemma EC.1, we get the desired result.  $\square$

We move onto the main result needed to prove Proposition EC.3. First, we build some terminology.

**DEFINITION EC.2 ( $\sigma$ -ADMISSIBLE LOSS FUNCTION).** A loss function  $\ell$  defined on  $\mathcal{Q} \times \mathcal{D}$  is  $\sigma$ -admissible with respect to  $\mathcal{Q}$  if the associated convex function  $c$  is convex in its first argument and the following condition holds:

$$\forall y_1, y_2 \in \mathcal{Y}, \forall d \in \mathcal{D}, |c(y_1, d) - c(y_2, d)| \leq \sigma |q_1 - q_2|,$$

where  $\mathcal{Y} = \{y : \exists q \in \mathcal{Q}, \exists \mathbf{x} \in \mathcal{X} : q(\mathbf{x}) = y\}$  is the domain of the first argument of  $c$ .

Note that  $\mathbb{R}^p$  is a reproducing kernel Hilbert space where the kernel is the standard inner product. Thus,  $\kappa$  in our case is  $X_{\max}$ .

*Proof of Proposition EC.3:* By the Lipschitz property of  $C(\cdot; d)$ ,

$$\sup_{d \in \mathcal{D}} |C(q_1(\mathbf{x}), d) - C(q_2(\mathbf{x}), d)| \leq (b \vee h) |q_1(\mathbf{x}) - q_2(\mathbf{x})|, \quad \forall q_1(\mathbf{x}), q_2(\mathbf{x}) \in \mathcal{Q}$$

as before, hence  $C : \mathcal{X} \times \mathcal{D} \rightarrow \mathbb{R}$  is  $(b \vee h)$ -admissible. Hence by Theorem 22 of Bousquet and Elisseeff (2002), the algorithm (NV-ERM2) has uniform stability with parameter  $\alpha_n^r$  as given.  $\square$

We have thus far established the stability of the machine learning algorithms (NV-ERM1) and (NV-ERM2), which lead to the risk bounds provided in Propositions EC.2 and EC.4, as follows.

**PROPOSITION EC.4 (Generalization Bound for (NV-ERM2)).** Let  $\hat{q}$  be the model produced by Algorithm (NV-ERM2) with  $\ell_2$  regularization. The following bound holds with probability at least  $1 - \delta$  over the random draw of the sample  $S_n$ , where each element of  $S_n$  is drawn iid from an unknown distribution on  $\mathcal{X} \times \mathcal{D}$ :

$$\frac{|R_{\text{true}}(\hat{q}) - \hat{R}(\hat{q}; S_n)|}{(b \vee h) \bar{D}} \leq \frac{(b \vee h)}{\bar{D} X_{\max}^{-2}} \frac{1}{n\lambda} + \left( \frac{2(b \vee h)}{\bar{D} X_{\max}^{-2}} \frac{1}{\lambda} + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}}. \quad (\text{EC.10})$$

We now state stability and generalization bound results for (NV-KO).

**PROPOSITION EC.5 (Uniform stability of (NV-KO)).** The algorithm (NV-KO) with iid data and the Gaussian kernel is symmetric with respect to the newsvendor cost function  $C(\cdot, \cdot)$  with uniform stability parameter

$$\alpha_\kappa = \frac{\bar{D}(b \vee h)^2}{(b \wedge h)} \frac{1}{1 + (n-1)r_w}, \quad (\text{EC.11})$$

where  $r_w = \exp(-2X_{\max}^2/w^2)$ .

The first term in the kernel stability parameter is the same as for (NV-ERM1), hence the insight that highly asymmetric backordering and holding costs leads to unstable decisions still holds. The second term depends on  $n$  and  $r_w$ , which in turn depends on the bandwidth  $w$  and the number of features  $p$  through  $X_{\max}$ . The term  $r_w$  goes from zero to one as  $w$  increases from zero to infinity;

this shows that greater stability is achieved with a larger bandwidth. This is an intuitive result, as a larger bandwidth is associated with bundling feature observations closer together. With this in mind we state the generalization bound for (NV-KO).

*Proof of Proposition EC.5:* The proof parallels that of Proposition EC.1. Symmetry follows from the fact that the data-generating process is i.i.d. For stability, we will change our notation slightly to make the dependence on  $n$  and  $S_n$  explicit. Let

$$q_n^\kappa = \arg \min_{q \geq 0} \tilde{R}(q; S_n, \mathbf{x}_{n+1}) = \arg \min_{q \geq 0} \frac{\sum_{j=1}^n \kappa_j [b(d_j - q)^+ + h(q - d_j)^+]}{\sum_{j=1}^n \kappa_j}$$

be the solution to (NV-KO) for the set  $S_n$ , and

$$q_{n \setminus i}^\kappa = \arg \min_{q \geq 0} \tilde{R}(q; S_n^{\setminus i}, \mathbf{x}_{n+1}) = \arg \min_{q \geq 0} \frac{\sum_{j \neq i} \kappa_j [b(d_j - q)^+ + h(q - d_j)^+]}{\sum_{j \neq i} \kappa_j}$$

be the solution to (NV-KO) for the set  $S_n^{\setminus i}$ . Note that:

$$\tilde{R}(q; S_n, \mathbf{x}_{n+1}) = \frac{\sum_{j \neq i} \kappa_j}{\sum_j \kappa_j} \tilde{R}(q; S_n^{\setminus i}, \mathbf{x}_{n+1}) + \frac{1}{\sum_j \kappa_j} \hat{R}(q; S_i, \mathbf{x}_{n+1}),$$

where  $S_i = (\mathbf{x}_i, d_i)$ .

By definition, the algorithm is stable if for all  $S_n \in \mathcal{Z}^n$  and  $i \in \{1, \dots, n\}$ ,

$$\sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} |C(q_n^\kappa, D(\mathbf{x})) - C(q_{n \setminus i}^\kappa, D(\mathbf{x}))| \leq \alpha_n,$$

where  $\alpha_n \leq O(1/n)$ . Now for a fixed  $\mathbf{x}$ , we have, by the Lipschitz property of  $C(q; \cdot)$ ,

$$\sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} |C(q_n^\kappa, D(\mathbf{x})) - C(q_{n \setminus i}^\kappa, D(\mathbf{x}))| \leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} (b \vee h) |q_n^\kappa - q_{n \setminus i}^\kappa|.$$

(See Fig. ??). So we want to bound  $|q_n^\kappa - q_{n \setminus i}^\kappa|$ .

By the convexity of the function  $\tilde{R}_n(\cdot; S_n, \mathbf{x}_{n+1})$ , we have (see Section 23 of Rockafellar (1997)):

$$\nu(q_{n \setminus i}^\kappa - q_n^\kappa) \leq \tilde{R}(q_{n \setminus i}^\kappa; S_n, \mathbf{x}_{n+1}) - \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})$$

for all  $\nu \in \partial \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})$  (set of subgradients of  $\tilde{R}(\cdot; S_n, \mathbf{x}_{n+1})$  at  $q_n^\kappa$ ). Further, because  $0 \in \partial \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})$  by the optimality of  $q_n^\kappa$ , we have

$$0 \leq \max_{\nu \in \partial \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})} \nu(q_{n \setminus i}^\kappa - q_n^\kappa) \leq \tilde{R}(q_{n \setminus i}^\kappa; S_n, \mathbf{x}_{n+1}) - \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})$$

where the max over  $\nu$  can be attained because  $\partial \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})$  is a compact set. Denote this maximum  $\nu^*$ .

Following arguments parallel those of the proof for (EC.1),

$$\begin{aligned} \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} (b \vee h) |q_n^\kappa - q_{n \setminus i}^\kappa| &\leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} \frac{b \vee h}{|\nu^*|} \cdot (\tilde{R}(q_{n \setminus i}^\kappa; S_n, \mathbf{x}_{n+1}) - \tilde{R}(q_n^\kappa; S_n, \mathbf{x}_{n+1})) \\ &\leq \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} \frac{(b \vee h) \kappa_i}{|\nu^*| \sum_j \kappa_j} \cdot |\tilde{R}(q_{n \setminus i}^\kappa; S_i, \mathbf{x}_{n+1})|. \end{aligned} \quad (\text{EC.12})$$

We can further simplify the upper bound (EC.12) as follows. Consider

$$\sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} \frac{\kappa_i}{\sum_j \kappa_j} = \sup_{(\mathbf{x}, D) \in \mathcal{X} \times \mathcal{D}} \frac{1}{1 + \sum_{j \neq i} \kappa_j / \kappa_i}.$$

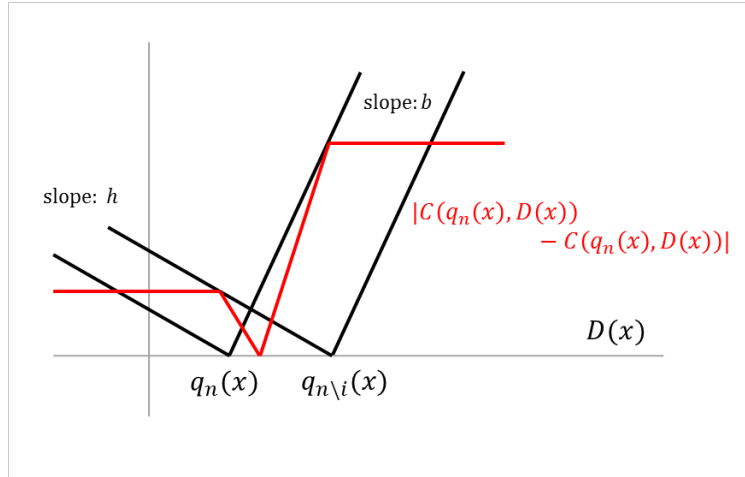
The supremum is thus achieved by the infimum of the ratio  $\kappa_j / \kappa_i$  over  $\mathcal{X} \times \mathcal{X} \times \mathcal{X}$ . For the Gaussian kernel,

$$\inf_{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_{n+1}) \in \mathcal{X} \times \mathcal{X} \times \mathcal{X}} \frac{K_w(\mathbf{x}_{n+1} - \mathbf{x}_j)}{K_w(\mathbf{x}_{n+1} - \mathbf{x}_i)} = \frac{e^{-4X_{\max}^2 p / 2w^2}}{e^0} = e^{-2X_{\max}^2 p / w^2} := r_w.$$

Finally, we can bound  $1/|\nu_j^*|$  by  $1/(b \wedge h)$ , as in the proof for Proposition EC.1. By using the tight uniform upper bound  $(b \vee h)\bar{D}$  on each term of  $|\tilde{R}(\cdot; \cdot, \mathbf{x}_{n+1})|$  from Lemma EC.1, we get the desired result.  $\square$

**PROPOSITION EC.6 (Generalization Bound for (NV-KO)).** *Let  $\hat{q}^\kappa$  be the optimal decision of (NV-KO) with the Gaussian kernel. The following bound holds with probability at least  $1 - \delta$  over the random draw of the sample  $S_n$ , where each element of  $S_n$  is drawn iid from an unknown distribution on  $\mathcal{X} \times \mathcal{D}$ :*

$$\frac{|R_{true}(\hat{q}^\kappa) - \hat{R}(\hat{q}^\kappa; S_n)|}{(b \vee h)\bar{D}} \leq \frac{2(b \vee h)}{b \wedge h} \frac{1}{1 + (n-1)r_w} + \left( \frac{4(b \vee h)}{b \wedge h} \frac{1}{1 + (n-1)r_w} + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}}.$$



**Figure EC.1** A plot illustrating that the difference  $|C(q_n(\mathbf{x}), D(\mathbf{x})) - C(q_{n \setminus i}(\mathbf{x}), D(\mathbf{x}))|$  is bounded.

We need the following lemma to wrap up the proofs of Propositions EC.2, EC.4 and EC.5.

**LEMMA EC.2.** *Let  $\hat{q}(S_n)$  be an in-sample decision with uniform stability  $\alpha_n$  with respect to a loss function  $\ell$  such that  $0 \leq \ell(\hat{q}(S_n), z) \leq M$ , for all  $z \in \mathcal{Z}$  and all sets  $S_n$  of size  $n$ . Then for any  $n \geq 1$  and any  $\delta \in (0, 1)$ , the following bound holds with probability at least  $1 - \delta$  over the random draw of the sample  $S_n$ :*

$$|R_{true}(\hat{q}(S_n)) - \hat{R}(\hat{q}(S_n), S_n)| \leq 2\alpha_n + (4n\alpha_n + M) \sqrt{\frac{\log(2/\delta)}{2n}}.$$

*Proof of Lemma EC.2:* The result is obtained by extending Theorem 12 of [Bousquet and Elisseeff \(2002\)](#) on pp. 507 by using the two-sided version of McDiarmid's inequality.  $\square$

We can now put the results together to prove Corollaries EC.2 and EC.4.

*Proof of Proposition EC.2:* The result follows from Proposition EC.1 and Lemma EC.2.  $\square$

*Proof of Proposition EC.4:* By Lemma EC.1,  $0 \leq \ell(A_S, z) \leq \bar{D}(b \vee h)$  for all  $z \in \mathcal{Z}$  and all sets  $S$ . The result then follows from Proposition EC.3 and Lemma EC.2.  $\square$

*Proof of Proposition EC.6:* The result follows from Proposition EC.5 and Lemma EC.2.  $\square$

Before proving Theorems 5, 6 and 7, we first prove the following lemma, that shows the equivalence of the newsvendor objective function to the nonparametric regression loss function by a multiplicative factor.

LEMMA EC.3. *The newsvendor objective function is a constant multiple of to the nonparametric regression loss function  $H_r(\cdot)$ .*

$$C(q; D) \equiv \frac{1}{2}H_r(q - D) := \frac{1}{2}[|q - D| + (2r - 1)(q - D)],$$

where  $r = b/(b + h)$ .

*Proof of Lemma EC.3:*

$$\begin{aligned} H_r(q - D) &:= |q - D| + (2r - 1)(q - D) \\ &= (q - D)^+ + (D - q)^+ + 2r(q - D) - (q - D) \\ &= (q - D)^+ + (D - q)^+ + 2r[(q - D)^+ - (D - q)^+] - (q - D)^+ + (D - q)^+ \\ &= 2r(q - D)^+ + 2(1 - r)(D - q)^+ \\ &= 2[r(q - D)^+ + (1 - r)(D - q)^+]. \quad \square \end{aligned}$$

*Proof of Theorem 5:*

$$\begin{aligned} |R_{true}(q^*) - \hat{R}_{in}(\hat{q}; S_n)| &\leq |R_{true}(\hat{q}) - \hat{R}_{in}(\hat{q}; S_n)| + |R_{true}(q^*) - R_{true}(\hat{q})| \\ &\leq \frac{2(b \vee h)^2 \bar{D}}{b \wedge h} \frac{p}{n} + \left( \frac{4(b \vee h)^2 \bar{D}}{b \wedge h} p + (b \vee h) \bar{D} \right) \sqrt{\frac{\log(2/\delta)}{2n}} + |R_{true}(q^*) - R_{true}(\hat{q})|, \end{aligned}$$

because the first term in the last inequality is the generalization bound in Proposition EC.2. As in the proof of Theorem 7, we can bound the second term by

$$|R_{true}(q^*) - R_{true}(\hat{q})| \leq (b \vee h) \mathbb{E}|q^* - \hat{q}|.$$

Finally, the term  $\mathbb{E}|q^* - \hat{q}|$  is the finite-sample bias of the decision  $\hat{q}$ , which is the  $b/(b + h)$  quantile of the conditional distribution of  $D|\mathbf{x}_{n+1}$  by Lemma EC.3. Then, by Theorem 3.2. of [Chaudhuri et al.](#)

(1991), we have  $|q^* - \hat{q}| = O(n^{m-k-\gamma/(2(k+\gamma)+p)})$  almost surely, where  $m$  is the order of derivatives (of the demand model) required to determine the true conditional quantile,  $k$  is the order of derivatives taken by the estimation procedure and  $0 < \gamma \leq 1$  is the Hölder continuity exponent of the demand model [see Chaudhuri et al. (1991) for the exact definitions]. For the linear demand model, we have  $\gamma = 1$ ,  $m = 0$  and  $k = 1$ , and we obtain the rate  $n^{-1/(2+p/2)}$ .

Computing the lower bound on the constant  $K$  is quite involved and below we present only the last steps in the computation. Following the steps in the proof of Theorem 3.2. of Chaudhuri et al. (1991), we require the convergence of the series

$$\sum_{n=1}^{\infty} n^{p+1-p(p+1)/(4+p)-c_7},$$

where  $c_7 = (1 - \varepsilon_1)^2 (c_5^*)^2 K^2 \tilde{c} / (p+1)$ , where  $\tilde{c} \in (c_1, c_2)$ ,  $c_1$  and  $c_2$  are as in the proof of Theorem 3.2. of Chaudhuri et al. (1991),  $c_5^*$  is as in the proof of Proposition 6.1. of Chaudhuri et al. (1991) and  $\varepsilon_2 \in (0, 1)$  is such that

$$\tilde{c} n^{4/(4+p)} - (p+1) \geq (1 - \varepsilon_2) \tilde{c} n^{4/(4+p)}$$

for all  $n > N_1$ , where  $N_1$  is required in Fact 6.5. of Chaudhuri et al. (1991). Rearranging and using the fact that over-harmonic series are convergent, we arrive at the requirement

$$K^2 > \frac{(8+5p)(1+p)}{4+p} \frac{1}{(1 - \varepsilon_1)^2 (c_5^*)^2 \tilde{c}}.$$

For a tighter bound, we want to find the smallest lower bound on  $K$ , i.e. the smallest possible lower bounds on  $\varepsilon_1$ , and the largest possible upper bounds on  $c_5^*$  and  $\tilde{c}$ .

From the proof of Proposition 6.1. of Chaudhuri et al. (1991), we have  $c_5^*$  equal to  $\lambda_2^* \rho / 2$ , where  $\rho$  is equal to  $1/3$  and  $\lambda_2^*$  is given by

$$\lambda_2^* := \min_{t \in [\underline{D}, \bar{D}]} f_\varepsilon(t),$$

which is bounded away from zero by assumption.

As for  $\varepsilon_1$ , we first rearrange the defining expression to get the requirement

$$\varepsilon_1 \geq \frac{p+1}{\tilde{c} n^{4/(4+p)}}.$$

From the proof of Proposition 6.1. of Chaudhuri et al. (1991),  $\tilde{c} n^{2/(4+p)} - (p+1)$  is a counting process so has to be greater than or equal to zero for all  $n \geq 1$ . Thus  $\tilde{c}$  is greater than or equal to  $(p+1)$ , and the lower bound for  $\varepsilon_1$  above can be upper bounded by  $2^{-4/(4+p)}$ , which is strict for  $n \geq 3$ , and so we can set  $\varepsilon_1^*$  equal to this value.

Combining, we have the lower bound

$$K \geq \sqrt{\frac{9(8+5p)}{(4+p)}} \frac{1}{(1 - 2^{-4/(4+p)}) \lambda_2^*},$$

and the best constant  $K$  is obtained by setting it equal to the right hand side.  $\square$



*Proof of Theorem 6:*

$$\begin{aligned}
|R_{true}(q^*) - \hat{R}_{in}(\hat{q}_\lambda; S_n)| &\leq |R_{true}(\hat{q}_\lambda) - \hat{R}_{in}(\hat{q}_\lambda; S_n)| + |R_{true}(\hat{q}) - R_{true}(\hat{q}_\lambda)| + |R_{true}(q^*) - R_{true}(\hat{q})| \\
&\leq (b \vee h) \left[ \frac{(b \vee h)X_{\max}^2 p}{n\lambda} + \left( \frac{2(b \vee h)X_{\max}^2 p}{\lambda} + \bar{D} \right) \sqrt{\frac{\log(2/\delta)}{2n}} \right] \\
&\quad + |R_{true}(\hat{q}) - R_{true}(\hat{q}_\lambda)| + |R_{true}(q^*) - R_{true}(\hat{q})|
\end{aligned}$$

because the first term in the last inequality is the generalization bound in Proposition EC.4. As in the proof of Theorem 5, we can bound the second and the third terms by

$$\begin{aligned}
|R_{true}(\hat{q}) - R_{true}(\hat{q}_\lambda)| &\leq (b \vee h) \mathbb{E}|\hat{q} - \hat{q}_\lambda|, \text{ and} \\
|R_{true}(q^*) - R_{true}(\hat{q})| &\leq (b \vee h) \mathbb{E}|q^* - \hat{q}|.
\end{aligned}$$

The first term  $\mathbb{E}|\hat{q} - \hat{q}_\lambda|$  is the finite-sample bias that results from regularization, which depends on the exact regularization used and on the demand distribution. The second term  $\mathbb{E}|q^* - \hat{q}|$  is the finite-sample bias of the decision  $\hat{q}$ , which we bound as in the proof for Theorem 5.  $\square$

*Proof of Theorem 7:* We have

$$\begin{aligned}
|R_{true}(q^*) - \hat{R}_{in}(\hat{q}^\kappa; S_n)| &\leq |R_{true}(\hat{q}^\kappa) - \hat{R}_{in}(\hat{q}^\kappa; S_n)| + |R_{true}(q^*) - R_{true}(\hat{q}^\kappa)| \\
&\leq \frac{2(b \vee h)^2 \bar{D}}{b \wedge h} \frac{1}{1 + (n-1)r_w} + \left( \frac{4(b \vee h)^2 \bar{D}}{b \wedge h} \frac{1}{1 + (n-1)r_w} + (b \vee h) \bar{D} \right) \sqrt{\frac{\log(2/\delta)}{2n}} \\
&\quad + |R_{true}(q^*) - R_{true}(\hat{q}^\kappa)|,
\end{aligned}$$

because the first term in the last inequality is the generalization bound in Proposition EC.6. As in the proof of Theorem 5, we can bound the second and the third terms by

$$\begin{aligned}
|R_{true}(\hat{q}) - R_{true}(\hat{q}^\kappa)| &\leq (b \vee h) \mathbb{E}|\hat{q} - \hat{q}^\kappa|, \text{ and} \\
|R_{true}(q^*) - R_{true}(\hat{q})| &\leq (b \vee h) \mathbb{E}|q^* - \hat{q}|.
\end{aligned}$$

The first term  $\mathbb{E}|\hat{q} - \hat{q}^\kappa|$  is the finite-sample bias that results from kernel-weights optimization, which depends on the choice of the kernel and on the demand distribution. The second term  $\mathbb{E}|q^* - \hat{q}|$  is the finite-sample bias of the decision  $\hat{q}$ , which we bound as in the proof for Theorem 5. The result then follows from similar arguments to the proof for Theorem 5.  $\square$

#### Appendix D: Retrieving sampling bound of [Levi et al. \(2007\)](#)

For the no-feature newsvendor problem ( $p = 1$  in our setup), Theorem 2.2. of [Levi et al. \(2007\)](#) presents a sampling size bound to obtain  $\epsilon$ -accurate decisions. Specifically, they state that when  $n \geq (9/2\epsilon^2)[(b+h)/(b \wedge h)]^2 \log(2/\delta)$ ,

$$R_{true}(\hat{q}) - R_{true}(q^*) \leq \epsilon R_{true}(q^*) \leq (b \vee h) \bar{D} \epsilon,$$

with probability at least  $1 - \delta$ , where  $0 < \delta < 1$  is a specified confidence level.

Let us now invert the generalization bound statement of Proposition EC.2 to arrive at a statement about the sampling size for  $\epsilon$ -accuracy. By the triangle inequality we have

$$|R_{true}(\hat{q}) - R_{true}(q^*)| \leq |R_{true}(\hat{q}) - R_{in}(\hat{q}; S_n)| + |R_{in}(\hat{q}; S_n) - R_{true}(q^*)|.$$

The first term is the generalization bound from Theorem EC.2; thus is less than

$$(b \vee h) \bar{D} \left[ \frac{2(b \vee h)}{(b \wedge h)n} + \left( \frac{4(b \vee h)}{b \wedge h} + 1 \right) \sqrt{\frac{\log(2/\delta)}{2n}} \right]$$

with probability at least  $1 - \delta$ . Inverting this statement, we have that

$$|R_{true}(\hat{q}) - R_{in}(\hat{q}; S_n)| \leq \epsilon_1 + \Delta_{1,n}$$

where  $\Delta_{1,n} = 2(b \vee h)^2 \bar{D} / (b \wedge h)n$ , whenever

$$n \geq [(b \vee h) \bar{D}] \left( \frac{4(b \vee h)}{b \wedge h} + 1 \right)^2 \frac{\log(2/\delta)}{2\epsilon_1^2}.$$

For the second term, consider, for some constant  $\epsilon_2 > 0$ ,

$$\begin{aligned} \mathbb{P}(|R_{in}(\hat{q}; S_n) - R_{true}(q^*)| \leq \epsilon_2) \\ &= \mathbb{P}(|R_{in}(\hat{q}; S_n) - R_{true}(q^*)| \leq \epsilon_2; |\hat{q} - q^*| \leq \Delta_{n,2}) + \mathbb{P}(|R_{in}(\hat{q}; S_n) - R_{true}(q^*)| \leq \epsilon_2; |\hat{q} - q^*| > \Delta_{n,2}) \\ &= \mathbb{P}(|R_{in}(\hat{q}; S_n) - R_{true}(q^*)| \leq \epsilon_2; |\hat{q} - q^*| \leq \Delta_{n,2}) \end{aligned}$$

where  $\Delta_{n,2} = Mn^{-2/(4+p)}\sqrt{\log n}$ , for some appropriately large  $M$  such that  $|\hat{q} - q^*| \leq \Delta_{n,2}$  almost surely (note such  $M$  exists because  $|\hat{q} - q^*| = O(n^{-2/(4+p)}\sqrt{\log n})$ ; see the proof of Theorem 5). We can further bound the above by

$$\begin{aligned} \mathbb{P}(|R_{in}(\hat{q}; S_n) - R_{true}(q^*)| \leq \epsilon_2; |\hat{q} - q^*| \leq \Delta_{n,2}) \\ &\geq \mathbb{P}(|R_{in}(\hat{q}; S_n) - R_{in}(q^*; S_n)| + |R_{in}(q^*; S_n) - R_{true}(q^*)| \leq \epsilon_2; |\hat{q} - q^*| \leq \Delta_{n,2}) \\ &\geq \mathbb{P}((b \vee h)|\hat{q} - q^*| + |R_{in}(q^*; S_n) - R_{true}(q^*)| \leq \epsilon_2; |\hat{q} - q^*| \leq \Delta_{n,2}) \\ &\geq \mathbb{P}(|R_{in}(q^*; S_n) - R_{true}(q^*)| \leq \epsilon_2 - (b \vee h)\Delta_{n,2}; |\hat{q} - q^*| \leq \Delta_{n,2}) \\ &\geq \mathbb{P}(|R_{in}(q^*; S_n) - R_{true}(q^*)| \leq \epsilon_2 - (b \vee h)\Delta_{n,2}). \end{aligned}$$

The final expression is a statement about the concentration of the sample average  $R_{in}(q^*; S_n)$  about its mean  $R_{true}(q^*)$ , so we can use the well-known Hoeffding's (concentration) inequality to get

$$|R_{in}(q^*; S_n) - R_{true}(q^*)| \leq \epsilon_2 - (b \vee h)\Delta_{n,2}$$

with probability at least  $1 - \delta$  whenever

$$n \geq (b \vee h) \bar{D} \frac{\log(2/\delta)}{2(\epsilon_2 - (b \vee h) \Delta_{n,2})^2}.$$

Combining, we have, with probability at least  $1 - \delta$ ,

$$|R_{true}(\hat{q}) - R_{true}(q^*)| \leq \epsilon_1 + \epsilon_2 + \Delta_{1,n} - (b \vee h) \Delta_{2,n}$$

whenever

$$n \geq (b \vee h) \bar{D} \max \left\{ \left( \frac{4(b \vee h)}{b \wedge h} + 1 \right)^2 \frac{\log(2/\delta)}{2\epsilon_1^2}, \frac{\log(2/\delta)}{2(\epsilon_2 - (b \vee h) \Delta_{n,2})^2} \right\}$$

Set  $\epsilon_1 = (b \vee h) \bar{D} \epsilon - \Delta_{1,n} = (b \vee h) \bar{D} [\epsilon - 4(b \vee h) / ((b \wedge h)n)]$  and  $\epsilon_2 = (b \vee h) \bar{D} \epsilon + (b \vee h) \Delta_{n,2}$ . Then we get

$$\frac{|R_{true}(\hat{q}) - R_{true}(q^*)|}{(b \vee h) \bar{D}} \leq \epsilon$$

whenever

$$\begin{aligned} n &\geq \frac{\log(2/\delta)}{2(b \vee h) \bar{D}} \max \left\{ \left( \frac{4(b \vee h)}{b \wedge h} + 1 \right)^2 \frac{1}{[\epsilon - 4(b \vee h) / ((b \wedge h)n)]^2}, \frac{1}{\epsilon^2} \right\} \\ &\geq \frac{1}{2(b \vee h) \bar{D}} \left( \frac{4(b \vee h)}{b \wedge h} + 1 \right)^2 \frac{\log(2/\delta)}{\epsilon^2} \end{aligned}$$

which is the same as the sampling bound of [Levi et al. \(2007\)](#) up to a constant factor (which bound is tighter depends on the exact values of  $b$ ,  $h$  and  $\bar{D}$ ).

## Appendix E: In-sample empirical results

No. of past days	without OS Features		with OS Features	
	Avg. Cost	Total no. of Features (avg. chosen)	Avg. Cost	Total no. of Features (avg. chosen)
0	0.9785	20 (3.0)	0.9785	20 (3.0)
1	0.9882	32 (4.0)	0.9848	32 (6.3)
2	0.9893	44 (5.3)	0.9966	44 (9.1)
3	0.9896	56 (6.2)	0.9335	56 (21.1)
4	0.9716	68 (8.0)	<b>0.8937</b>	68 (28.1)
5	0.9711	80 (8.3)	0.9112	80 (36.7)
6	0.9700	92 (8.6)	0.9080	92 (43.4)
7	0.9691	104 (9.5)	0.9270	104 (52.4)
8	0.9688	116 (9.5)	0.9097	116 (56.5)
9	0.9687	128 (9.6)	0.9329	128 (62.6)
10	0.9690	140 (9.6)	0.9195	140 (69.7)
11	0.9693	152 (9.7)	0.9204	152 (73.4)
12	<b>0.9685</b>	164 (9.9)	0.9459	164 (82.2)
13	0.9689	176 (10.3)	0.8976	176 (91.8)
14	0.9686	188 (10.3)	0.9002	188 (97.2)

**Table EC.1** Average in-sample cost of the solution to (NV-ERM1) with the day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks), with and without OS features. The best in-sample results are highlighted in bold. The total number of features that appear in the decision are also reported (as measured by the number of nonzero elements in the decision vector).

No. of clusters	Avg. Cost
2	<b>1.3167</b>
3	1.3727
4	1.4332
5	1.3977
6	1.4170
7	1.4244
8	1.4439
9	1.4455
10	1.4293
11	1.4583
12	1.4343

**Table EC.2** Average in-sample cost of the Cluster + SAA approach, with OS features. The best in-sample results are highlighted in bold. We include all OS features as not including them often yields empty clusters.

No. of past days	without OS Features		with OS Features	
	Avg. Cost	Total no. of Features (avg. chosen)	Avg. Cost	Total no. of Features (avg. chosen)
0	1.0945	20 (16.7)	1.0945	20 (16.7)
1	<b>1.0892</b>	32 (16.8)	1.0381	44 (27.3)
2	1.0960	44 (16.6)	1.0775	68 (37.6)
3	1.1142	56 (16.6)	0.9372	92 (48.5)
4	1.1220	68 (16.8)	0.9393	116 (59.5)
5	1.1360	80 (16.4)	0.9270	140 (71.2)
6	1.1525	92 (16.3)	<b>0.9202</b>	164 (83.1)
7	1.1487	104 (16.5)	0.9632	188 (95.4)
8	1.1596	116 (16.4)	1.0301	212 (108.3)
9	1.1638	128 (16.6)	1.0500	236 (118.6)
10	1.1587	140 (16.5)	1.0521	260 (129.7)
11	1.1694	152 (16.4)	1.0899	284 (141.4)
12	1.1601	164 (16.5)	1.0799	308 (153.5)
13	1.1540	176 (27.7)	1.0464	332 (165.5)
14	1.1658	188 (30.7)	1.0785	356 (178.2)

**Table EC.3** Average in-sample cost of the solution to the SEO approach with day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks), with and without OS features. The best in-sample results are highlighted in bold. The total number of features that appear in the decision are also reported (as measured by the number of elements in the decision vector that are at least as large as 1% of the largest element in absolute value).

Regularization Param.	without OS Features		with OS Features	
	Avg Cost	Total no. of Features (avg. chosen)	Avg Cost	Total no. of Features (avg. chosen)
$1 \times 10^{-4}$	1.1450	188 (9.1)	0.6140	356 (4.7)
$5 \times 10^{-5}$	1.0769	188 (12.0)	0.5793	356 (5.8)
$1 \times 10^{-5}$	1.0750	188 (14.6)	0.5837	356 (7.2)
$5 \times 10^{-6}$	1.0748	188 (15.3)	0.5837	356 (7.8)
$1 \times 10^{-6}$	1.0117	188 (14.6)	0.5376	356 (10.3)
$5 \times 10^{-7}$	0.9737	188 (13.0)	0.5165	356 (11.5)
$1 \times 10^{-7}$	<b>0.9729</b>	188 (13.2)	<b>0.4489</b>	356 (28.1)

**Table EC.4** Average in-sample cost of the solution to (NV-ERM2) solved with  $\ell_1$  regularization, with the day of the week and time of the day features and 2 weeks of past demands with and without OS features for a range of regularization parameters. The best in-sample results are highlighted in bold. The total number of features that appear in the decision are also reported (as measured by the number of nonzero elements in the decision vector).

Regularization Param.	without OS Features		with OS Features	
	Avg. Cost	Total no. of Features (avg. chosen)	Avg. Cost	Total no. of Features (avg. chosen)
$1 \times 10^{-4}$	1.1138	188 (8.7)	1.1153	356 (8.7)
$5 \times 10^{-5}$	1.0626	188 (11.2)	1.0629	356 (11.2)
$1 \times 10^{-5}$	1.0462	188 (13.0)	1.0463	356 (13.0)
$5 \times 10^{-6}$	1.0434	188 (13.4)	1.0438	356 (13.7)
$1 \times 10^{-6}$	0.9805	188 (13.2)	0.9276	356 (17.7)
$5 \times 10^{-7}$	<b>0.9570</b>	188 (10.1)	0.9270	356 (21.6)
$1 \times 10^{-7}$	0.9684	188 (10.2)	<b>0.9153</b>	356 (42.5)

**Table EC.5** Average in-sample cost of the solution to (NV-ERM2) solved with  $\ell_2$  regularization, with the day of the week and time of the day features and 2 weeks of past demands with and without OS features for a range of regularization parameters. The best in-sample results are highlighted in bold. The total number of features that appear in the decision are also reported (as measured by the number of elements in the decision vector that are at least as large as 1% of the largest element in absolute value).

Regularization Param.	without OS Features		with OS Features	
	Avg Cost	Total no. of Features (avg. chosen)	Avg Cost	Total no. of Features (avg. chosen)
$1 \times 10^0$	1.1020	188 (9.9)	1.4849	356 (9.9)
$5 \times 10^{-1}$	<b>1.0905</b>	188 (10.4)	1.2418	356 (10.4)
$1 \times 10^{-1}$	1.0990	188 (14.9)	1.2318	356 (14.9)
$5 \times 10^{-2}$	1.1023	188 (19.7)	1.1975	356 (15.9)
$1 \times 10^{-2}$	1.1264	188 (34.0)	1.1307	356 (24.1)
$5 \times 10^{-3}$	1.1365	188 (44.2)	<b>1.0636</b>	356 (28.5)

**Table EC.6** Average in-sample cost of the solution to the SEO approach solved with  $\ell_1$  regularization, with day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks), with and without OS features. The best in-sample results are highlighted in bold. The total number of features that appear in the decision are also reported (as measured by the number of elements in the decision vector that are at least as large as 1% of the largest element in absolute value).

Regularization Param.	without OS Features		with OS Features	
	Avg Cost	Total no. of Features (avg. chosen)	Avg Cost	Total no. of Features (avg. chosen)
$1 \times 10^{-1}$	<b>1.1486</b>	188 (57.6)	1.1584	356 (92.8)
$5 \times 10^{-2}$	1.1555	188 (61.2)	1.1357	356 (85.8)
$1 \times 10^{-2}$	1.1613	188 (65.0)	1.0398	356 (88.0)
$5 \times 10^{-3}$	1.1626	188 (65.7)	<b>1.0040</b>	356 (87.1)
$1 \times 10^{-3}$	1.1669	188 (67.8)	1.0526	356 (76.2)
$5 \times 10^{-4}$	1.1658	188 (68.6)	1.0636	356 (69.9)
$1 \times 10^{-4}$	1.1654	188 (69.2)	1.0716	356 (64.3)

**Table EC.7** Average in-sample cost of the solution to the SEO approach solved with  $\ell_2$  regularization, with day of the week and time of the day features and an increasing number of days of past demands (for up to 2 weeks), with and without OS features. The best in-sample results are highlighted in bold. The total number of features that appear in the decision are also reported (as measured by the number of elements in the decision vector that are at least as large as 1% of the largest element in absolute value).