

SHAPE-CONSTRAINED KERNEL PDF AND PMF ESTIMATION

Pang Du, Christopher F. Parmeter, Jeffrey S. Racine*

Virginia Tech, University of Miami and McMaster University

Abstract: We present an approach for estimating shape-constrained kernel-based probability density functions (PDFs) and probability mass functions (PMFs) that includes constraints on the PDF (PMF) function itself, its integral (sum), and derivatives (finite differences) of any order. We also allow for pointwise upper and lower bounds (i.e., inequality constraints) on the PDF and PMF, in addition to more popular equality constraints. Furthermore the approach handles a range of transformations of the PDFs and PMFs including, for example, logarithmic transformations, which allow us to impose log-concave or log-convex constraints. We also provide the theoretical underpinnings for the procedures. The results of a simulation-based comparison between our proposed approach and those Grenander-type methods favor our approach when the data-generating process is smooth. To the best of our knowledge, ours is also the only *smooth* framework that handles PDFs and PMFs in the presence of inequality bounds, equality constraints, and other popular constraints. An implementation in R incorporates constraints such as monotonicity (both increasing and decreasing), convexity and concavity, and log-convexity and log-concavity, among others, while respecting finite-support boundaries by using boundary kernel functions.

Key words and phrases: Kernel density estimation, probability density function, probability mass function, shape constraints.

1. Introduction

Shape constraints play a vital role in identification, estimation, and inference in econometric and statistical applications (see, e.g., Chetverikov, Santos and Shaikh (2018) for a review of recent developments and their importance in applied work). Such constraints sometimes emerge naturally owing to the nature of the data, but increasingly often are required when replacing parametric models with more versatile semi- and nonparametric models. The ability to preserve the qualitative shape properties present in a parametric model is a key component of any alternative method. However, the consequences of misspecifying the parametric model can be severe, and influence the choice of the nonparametric alternative. There are two reasons why one might wish to integrate shape constraints into a nonparametric estimation procedure. The first is to achieve potential gains in estimator *efficiency* by imposing *valid* shape constraints on

*Corresponding author.

some statistical object of interest. That is, if one's assumption about a shape constraint on an otherwise unspecified curve is correct, then incorporating this information into the estimation procedure can improve the finite-sample performance of the corresponding estimator. The second reason is to assess the validity of the shape constraints using formal quantitative inference, or to determine the qualitative effect of the constraints on the resulting estimate.

Imposing shape constraints on an otherwise unrestricted nonparametric curve is a key element of a sound empirical analysis that encompasses a range of approaches; see Groeneboom and Jongbloed (2014) for examples of shape-constrained estimators and algorithms, along with their theoretical properties. Perhaps the most common applications of enforcing shape constraints arise when modeling a conditional mean function (i.e., a regression), which is understandable, given the popularity of regression analysis. However, the density function is also a popular object of interest that necessitates a separate treatment from that of regression, owing to its unique nature. Shape-constrained density estimation, like its regression-based counterpart, has a rich history that can be traced to the seminal work of Grenander (1956), who analyzes the maximum likelihood estimator (MLE) of a decreasing density on the nonnegative half-line (see also Groeneboom and Jongbloed (2018) for recent theoretical work in this direction). Note that Prakasa Rao (1969) shows that this estimator exhibits nonstandard asymptotic behavior, because it converges at a cube rate ($n^{-1/3}$) at points at which the true decreasing density is differentiable with a negative derivative. This is slower than competing local kernel-based estimators that assume smoothness ($n^{-2/5}$), a common assumption among practitioners that we adopt in one of the two kernel-based estimators we consider here. Although the density function is our main object of interest, we also treat the mass function, and note that kernel-based mass function estimators for categorical data have a different (and faster, i.e., $n^{-1/2}$) rate of convergence than that of their kernel density-based counterparts.

In density estimation settings, a variety of innovative approaches have been proposed for imposing specific constraints, such as monotonicity, concavity, and log-concavity, among others. Though some of these approaches admit certain combinations of shape constraints, many are tailored to a *particular* setting (e.g., monotonicity *only*). In addition, while some existing approaches incorporate bounds on the *support* of the variable under study, others do not. Furthermore, most existing approaches are predicated on *continuously* distributed random variables, though constrained probability mass functions (PMFs) may also be of value when modeling *discrete support* random variables, which arise frequently in applied settings.

Grenander-based approaches (Grenander (1956)) have been widely used to impose certain shape constraints, and one of their appealing features is that they do not require any *tuning parameters*, unlike *smooth* kernel-based nonparametric

methods, such as those proposed below, which require the specification of a *bandwidth*. However, although Grenander-based approaches are nonparametric in nature, they are *nonsmooth* which runs counter to the spirit of adopting a *smooth* nonparametric approach in the first place. For example, the approach Grenander (1956) proposes for imposing monotonicity can be characterized as the left derivative of the least concave majorant of the empirical distribution function, which is a nonsmooth function. Practitioners who routinely assume smoothness and adopt smooth nonparametric estimators are not likely to be attracted to nonsmooth nonparametric shape-constrained solutions, hence the appeal of *smooth* shape-constrained nonparametric solutions, such as those proposed herein.

The literature on constrained nonparametric estimation has grown significantly over the past few decades. The approach that we extend here has proven to be a particularly popular, versatile, and extensible method for imposing constraints on a *smooth* nonparametric object (see Hall and Presnell (1999)). This approach places weights directly on the sample realizations so that the desired constraint is imposed effectively. In kernel-based *regression* settings, this amounts to starting with a standard kernel estimator. Then, if the constraints are violated in some region of the support we shift the regressand *vertically* in such a way that a standard kernel regression on the *shifted* regressand delivers a regression curve that satisfies the required constraints, while minimizing some distance metric from the unconstrained regression function (Hall et al. (2001); Du, Parmeter and Racine (2013)). In kernel-based *density* settings, this approach can be leveraged by placing weights on the *kernel function* associated with each sample realization (as opposed to the sample realizations themselves) to produce a density that satisfies the required constraints. A similar method, known as *data sharpening* (Hall and Kang (2005)), instead introduces weights that shift the data *horizontally* prior to smoothing, a subtle, but important distinction. We adopt the approach of Hall and Presnell (1999), because vertically shifting observations can be undertaken using standard off-the-shelf quadratic programming methods, whereas horizontally shifting observations may require full-blown nonlinear programming, which may be less tractable from a practical perspective.

Building on the work of Du, Parmeter and Racine (2013), who consider a unified framework for *smooth* shape-constrained nonparametric kernel regression, we propose a unified framework for *smooth* shape-constrained kernel density and PMF estimation. Shape-constrained kernel density (and mass) function estimation differs from shape-constrained kernel regression in terms of both its practical implementation and in its theoretical properties, and hence requires a separate treatment. Our approach is extremely flexible, and allows for a range of constraints to be imposed *simultaneously* (presuming, of course, that the set of constraints is internally consistent). The original implementation

(Hall et al. (2001)) involves optimizing a power-divergence criterion. Du, Parmeter and Racine (2013) propose replacing this criterion with an L_2 -norm criterion, which delivers an estimator that retains all of the desirable features of the power-divergence-based method, but is far more flexible and extensible and far simpler to solve from a practical perspective. The method proposed here generalizes the seminal work of Hall and Huang (2002), who impose unimodality on a univariate kernel density estimator, and modify it in such a way as to deliver a unified approach with a straightforward implementation. We believe that this unified framework will be of particular interest to practitioners who wish to simultaneously impose a range of constraints in a smooth nonparametric setting.

Additionally, we build on the insights of Li, Liu and Li (2017), who propose a slightly modified version of the optimization criterion proposed by Du, Parmeter and Racine (2013). While Li, Liu and Li (2017) adopt an L_2 -norm criterion, as per Du, Parmeter and Racine (2013), rather than optimizing the distance between the optimization *weights* and their unconstrained counterparts, they instead optimize the distance between the constrained *estimates* and their unconstrained counterparts. Although Li, Liu and Li (2017) provide convincing simulation evidence that their modification can deliver constrained estimates with improved finite-sample performance, they offer no theoretical justification for this modification. We demonstrate theoretically that this modified L_2 -optimization criterion delivers constraint weights that ensure *identical* asymptotic behavior to that from optimizing the weights directly. By providing the theoretical underpinnings for the slightly modified optimization criterion proposed by Li, Liu and Li (2017), we establish that the constraint weights can be based on this criterion with no loss of information.

Finally, we demonstrate how our method can be adapted to handle constraints on the *log-density*. This is an important generalization, because constraints on the log-density, when enforced using the density function directly, can result in a difficult nonlinear optimization problem. By focusing instead *directly* on the log-density, we ensure straightforward constraint enforcement, with trivial conversion back to the constrained density itself, all within the same unified theoretical framework as that for constraints directly on the density.

The proposed approach differs from that of Du, Parmeter and Racine (2013), among others, in several ways. In our setting, we are dealing with density estimation and weights are applied on the kernel function. In contrast, in the regression setting of Du, Parmeter and Racine (2013), weights are applied on the dependent variable, which affects the proofs in a nontrivial way. Here we prove Theorem 2 for the Cramér–von Mises distance function (earlier works have not considered this distance metric), which requires handling cross-product terms involving the constraint weights in the various components of our decomposition of the constrained density estimator. Additionally, Theorem 3 is entirely new. To

the best of our knowledge, it represents the first attempt to impose smoothness constraints on a PMF estimator. While not a theoretical contribution, we also demonstrate how to impose log-concavity on a smooth kernel density estimate in a simple quadratic programming setup.

One of the constraints on the log-density, specifically *log-concavity*, has long been a topic of interest in statistics; see Walther (2009) for an introduction, and Samworth and Sen (2018) for a recent review. Briefly, log-concave densities present an appealing and natural alternative to the class of unimodal densities. Though the class of log-concave densities is a subset of the class of unimodal densities, it contains most of the commonly used parametric distributions, and is therefore a rich and useful nonparametric class. Recent developments include the works of Feng et al. (2021) who study an adaptation of the nonparametric MLE density for the class of upper semi-continuous log-concave densities on \mathbb{R}^d (the logarithm of the resulting estimate is a *piecewise-linear* nonsmooth function), and Rathke and Schnörr (2019), who propose a fast implementation of the smoothed version of this estimator.

Log-concavity has also played an important role in applied microeconomic analysis. By imposing log-concavity in an otherwise unrestricted nonparametric setting, economic studies that previously relied on a specific parametric model can instead rely on less restrictive nonparametric models leading to more robust results. Examples include the works of Bagnoli and Bergstrom (2005), who describe how the log-concavity assumption allows *just enough* special structure to yield workable theories across various subfields, Meyer-ter-Vehn, Smith and Bognar (2017), who explore costly deliberations by two differentially informed and possibly biased jurors, exploiting an assumption that jurors' information types have a log-concave density, and Tan and Zhou (2020), who rely on log-concavity in agent heterogeneity to establish several formal results in a model of price competition entry and multi-sided markets.

Our adaptation of the work of Hall and Huang (2002) to log-concavity also stands in contrast to a recently proposed kernel-based linear adjustment mechanism (Wolters and Braun (2018a,b)) that tackles constrained estimation using a specified number of inflection points. This approach can also be used to enforce log-concavity, though the authors do not consider this particular constraint. However, it would require that we know the locations of these inflection points *ex ante*, otherwise they need to be approximated using some optimization routine, which has its drawbacks. In contrast, our proposed approach to imposing log-concavity requires no prior knowledge or approximations of the locations of the inflection points. Instead, we impose the constraints on the log-density directly, leading to a direct system of linear inequality constraints, and thus a fast and efficient algorithm for imposing log-concavity in a smooth setting is provided. The linear adjustment mechanism of Wolters and Braun (2018b) is equivalent to our approach *if* the number of adjustment functions is equivalent to the number

of observations *and* the adjustment functions themselves are equivalent to the kernel smoothing function of the unconstrained density estimator. However, they do not consider using their linear adjustment mechanism approach to impose log-concavity which, given its popularity in applied settings, forms the basis for one of the Monte Carlo simulations we run to compare our approach with its peers; see the R package `scdensity` (Wolters (2018)) for an implementation of the linear adjustment mechanism approach.

In addition to the works referenced above, the related literature includes the studies of Woodroffe and Sun (2002), who consider a penalized MLE estimate of a density on the positive half of the real number line when the density is nonincreasing, Meyer and Woodroffe (2004), who develop a nonparametric MLE that is consistent for the mode, Hall and Kang (2005), who consider unimodal kernel density estimation using data sharpening, Dette and Pilz (2006), who conduct a comparative study of monotone constrained estimators, Birke (2009), who considers shape-constrained density estimation using monotone rearrangement (Hardy, Littlewood and Pólya (1952); Chernozhukov, Fernandez-Val and Galichon (2009)), Dümbgen and Rufibach (2009), who studied the MLE of a log-concave density, and Koenker and Mizera (2010), who formulate the MLE of a log-concave density as a convex optimization problem, showing that it has an equivalent dual formulation to that of a constrained maximum Shannon entropy problem. Cule, Samworth and Stewart (2010) study a nonsmooth log-concave MLE of a probability distribution function and Meyer and Habtzghi (2011) use regression splines, based on the work of Meyer (2008), to formulate a nonparametric MLE of strictly decreasing probability densities in terms of convex programming and iteratively re-weighted least squares cone projection algorithms. Chen and Samworth (2013) study the smoothed log-concave MLE of a probability distribution function, Horowitz and Lee (2017) explain how to estimate and obtain an asymptotic uniform confidence band for a conditional mean function under possibly nonlinear shape restrictions, and Koenker and Mizera (2018) consider a log-concave estimation for weaker forms of concavity constraints that allow for heavier tail behavior and sharper modal peaks. More recently, Lok and Tabri (2002) develop an empirical tilting method for shape-constrained estimation over a data-driven grid of points to enforce the stochastic dominance of a pair of cumulative distribution functions.

The rest of this paper proceeds as follows. Section 2 presents a unified framework for kernel-based probability density function (PDF) and PMF estimators and describes our approach. Here, Section 2.1 examines how we determine the constraint weights, and Section 2.2 briefly discusses finite-support boundary kernel functions and presents several examples of popular constraints. Section 3 outlines the theoretical properties of the proposed approach and Section 4 presents a set of Monte Carlo simulations that show that the proposed approach is competitive with, and often improves upon leading methods that have been

tailored to two popular constraints (log-concavity and monotonicity). Section 5 concludes the paper. Detailed theoretical proofs are relegated to the online Supplementary Material, and an open implementation in R exists to assist practitioners interested in exploring the proposed methods.

2. Shape-Constrained Kernel Density Estimation

Let X_i , for $i = 1, \dots, n$ be an independent and identically distributed (i.i.d.) random sample drawn from $f(x)$, where n denotes the sample size. To estimate $f(x)$ using smooth nonparametric methods, we begin with the standard kernel density estimator,

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right), \quad (2.1)$$

where h is the *bandwidth*, $K(\cdot)$ is the *kernel function*, usually chosen as a symmetric mean zero PDF itself, and x is a support point at which the density is estimated (Rosenblatt (1956); Parzen (1962)). To help discuss our (constraint) weighted density estimator, when imposing constraints on the density function, we introduce a vector of constraint weights p_i , for $i = 1, \dots, n$, and modify (2.1) as follows:

$$\hat{f}(x|p) = \frac{1}{h} \sum_{i=1}^n p_i K\left(\frac{x - X_i}{h}\right). \quad (2.2)$$

Note that for $p_i = p_{unif} = 1/n$, the *uniform* weights $\hat{f}(x|p_{unif}) = \hat{f}(x)$, which is the standard (i.e., *unconstrained*) estimator (2.1). In other words, we use the notation $\hat{f}(x|p_{unif})$ in what follows to represent (2.2) for the special case in which the constraint weights assume the value $p_i = 1/n$, for $i = 1, \dots, n$. Furthermore, these special weights are denoted as p_{unif} , and for these *and only these*, weights (2.2) is equal to (2.1), the standard kernel estimator (which we call the unconstrained estimator).

To impose constraints on the density function, we let $p_i = n^{-1}(1 + a_i)$ act as the constraint weights in (2.2), yielding the estimator

$$\begin{aligned} \hat{f}(x|p) &= \frac{1}{nh} \sum_{i=1}^n (1 + a_i) K\left(\frac{x - X_i}{h}\right) \\ &= \frac{1}{nh} \sum_{i=1}^n (1 + a_i) K(Z_i) \\ &= \frac{1}{nh} \sum_{i=1}^n K(Z_i) + \frac{1}{nh} \sum_{i=1}^n a_i K(Z_i) \\ &= \hat{f}(x|p_{unif}) + \frac{1}{nh} \sum_{i=1}^n a_i K(Z_i), \end{aligned} \quad (2.3)$$

where $Z_i = (x - X_i)/h$ and the *unconstrained* (i.e., *uniform*) weights are $a_i = 0$ (i.e., $p_i = 1/n$, the weights that return the unconstrained estimator).

Imposing constraints on the log-density function can be accomplished with a slightly modified setup. To impose constraints on the log-density function (or its derivatives), we instead consider an estimator of the form

$$\hat{f}(x|p) = \hat{f}(x) \prod_{i=1}^n \exp \left\{ \frac{a_i K(Z_i)}{nh} \right\}. \quad (2.4)$$

Taking the logarithm of both sides, we obtain

$$\log \hat{f}(x|p) = \log \hat{f}(x) + \frac{1}{nh} \sum_{i=1}^n a_i K(Z_i).$$

Hence the constrained density estimator when imposing constraints on the log-density is given by

$$\hat{f}(x|p) = \exp \left\{ \log \hat{f}(x|p_{unif}) + \frac{1}{nh} \sum_{i=1}^n a_i K(Z_i) \right\},$$

where, the *unconstrained* weights used in the object $\hat{f}(x|p_{unif})$, correspond to $a_i = 0$ in (2.4) which delivers the standard kernel density estimator $\hat{f}(x)$. Regardless of the constraints considered, any constraints imposed on either the density or the log-density, as expressed above, will be *linear* in a_i , which, combined with a quadratic objective function, leads naturally to solving a quadratic program. The resulting constrained estimator arises from solving a quadratic program, and then replacing the arbitrary weights a_i with the feasible constrained weights determined by the quadratic program.

Thus far, we have outlined two approaches that introduce weights that deliver constrained density or log-density estimates. Now, we explicitly introduce the constraints themselves in a general framework. Denote the j th derivative of $\hat{f}(x|p)$, $\log \hat{f}(x|p)$, and $K(Z_i)$ with respect to x as $\hat{f}^{(j)}(x|p)$, $\log^{(j)} \hat{f}(x|p)$, and $K^{(j)}(Z_i)$, respectively (the same goes for $\hat{f}(x|p_{unif})$ and $\log \hat{f}(x|p_{unif})$). Let $l(x)$ and $u(x)$ denote *pointwise* lower and upper bounds, respectively, that may change with x , where $l(x) \leq u(x)$. The constraints on the j th derivative of the density and log-density, for $j = 0, 1, 2, \dots$, can be expressed as

$$l(x) \leq \hat{f}^{(j)}(x|p) \leq u(x) \quad (2.5)$$

and

$$l(x) \leq \log^{(j)} \hat{f}(x|p) \leq u(x),$$

respectively. Thus, for $j = 0$, we are constraining the density or log-density, for $j = 1$, we are constraining the first derivative thereof, and so on. Consider, by

way of illustration, the constraint $\hat{f}^{(j)}(x|p) \geq l(x)$, which we express as

$$\hat{f}^{(j)}(x|p_{unif}) + \frac{1}{nh} \sum_{i=1}^n a_i K^{(j)}(Z_i) \geq l(x)$$

or

$$\frac{1}{nh} \sum_{i=1}^n a_i K^{(j)}(Z_i) \geq l(x) - \hat{f}^{(j)}(x|p_{unif}).$$

Furthermore, the constraint $\log^{(j)} \hat{f}(x|p) \geq l(x)$ (the lower bound $l(x)$ may well differ from that for $\hat{f}^{(j)}(x|p)$ above) can be expressed as

$$\frac{1}{nh} \sum_{i=1}^n a_i K^{(j)}(Z_i) \geq l(x) - \log^{(j)} \hat{f}(x|p_{unif}).$$

One appealing feature of our approach is that we can *simultaneously* impose a set of internally consistent constraints. For instance, if we wish to impose the constraints that $\hat{f}^{(0)}(x|p) = \hat{f}(x|p) \geq 0$ (nonnegativity of the constrained density) and $\log^{(2)} \hat{f}(x|p) \leq 0$ (log-concavity), we can impose the constraints

$$\frac{1}{nh} \sum_{i=1}^n a_i K(Z_i) \geq -\hat{f}(x|p_{unif})$$

and

$$-\frac{1}{nh} \sum_{i=1}^n a_i K^{(2)}(Z_i) \geq \log^{(2)} \hat{f}(x|p_{unif}).$$

When solving the quadratic program outlined in the next section, we typically impose the constraint $\sum_{i=1}^n a_i = 0$.

We wish to handle a rich array of constraints, and we may also find ourselves in settings with random variables having either unbounded or compact support. The most popular approaches for compact support kernel estimation use one kernel function when the support is bounded above and below (e.g., Beta(a,b)), one when the support is bounded below (e.g., Gamma(a)), or multiple kernel functions when the support is bounded above and below (e.g., floating boundary kernel functions). To deal with compact support random variables, in Section 2.2, we use *kernel carpentry* to provide a flexible kernel function that is well-suited to the current setting.

2.1. Selection of the constraint weights

Having established how to construct the constrained estimator for an *arbitrary* set of weights, we now examine how best to select the weights to satisfy some *particular* constraint of interest.

A variety of approaches for constrained weight selection have been proposed in the literature, each of which minimizes some measure of *divergence* between the constrained and the unconstrained *weights* or the constrained and the unconstrained *estimates*. Some divergence metrics are more computationally demanding than others, and different metrics may impose binding restrictions on the weights in order to produce valid estimates. For example, Hall et al. (2001) suggest using the Cressie–Read power-divergence metric, Hall and Huang (2002) investigate a smoothed Cramér–von Mises metric, and Du, Parmeter and Racine (2013) suggest an L_2 -norm metric. Specifically, in the power-divergence and L_2 -norm frameworks, the constrained weights are selected to be as close as possible to the unconstrained weights (also called the *uniform* weights), whereas in the smoothed Cramér–von Mises setting, the constrained weights are chosen to minimize the squared integrated difference between the unconstrained and the constrained densities. As Hall and Huang (2002) and Du, Parmeter and Racine (2013) document, one benefit of adopting an L_2 -norm (i.e., the squared distance) metric is that we can select smoothing parameters based solely on the unconstrained estimator. Hence, standard off-the-shelf methods can be used without modification, and we maintain this practice in what follows.

Following Hall and Huang (2002), Du, Parmeter and Racine (2013), and Li, Liu and Li (2017), we consider two closely related approaches for the optimal construction of the constraint weights, and emphasize their relative strengths. For the first approach, we minimize the L_2 -norm divergence between the constraint weights and the uniform weights, where the divergence metric is defined as follows:

$$D_{L_2}(p) = (p_u - p)'(p_u - p).$$

In this case, provided the desired constraints are *linear* in p , we can solve this minimization problem by means of a straightforward quadratic program exercise using, say, the *quadprog* package (Turlach and Fortran (2019)) in R. For the second approach, we minimize a smoothed Cramér–von Mises distance metric, where the squared integrated difference between the unconstrained and constrained densities is defined as follows:

$$D_{CM}(p) = (n^2|h|)^{-1} \sum_{i=1}^n \sum_{j=1}^n (np_i - 1)(np_j - 1)L\left(\frac{X_i - X_j}{h}\right), \quad (2.6)$$

where $L(\cdot)$ is the convolution kernel of $K(\cdot)$ with itself. Regardless of the metric used, $D_{L_2}(p)$ and $D_{CM}(p)$ require that the constraint weights themselves satisfy a constraint in order to guarantee that a proper probability density is produced (i.e., for constraints on the density function using (2.3) or constraints on the log-density function using (2.4), we require $\sum_{i=1}^n a_i = 0$). It is useful to focus on the relative merits of each distance metric used to select the constraint weights. The obvious benefit of $D_{L_2}(p)$ is the relative theoretical ease with which to assess the properties

of the corresponding constraint weights. As Du, Parmeter and Racine (2013) demonstrate, the relative magnitude of the constraint weights with the L_2 -norm is $O(n^{-1})$. We can also view a difference from the uniform weights as a measure of relative entropy with respect to the uniform distribution. The Cramér-von Mises metric has obvious practical appeal, because it selects weights that lead to the constrained density deviating as little as possible from the unconstrained density. Moreover, as shown in simulations here and in Li, Liu and Li (2017), selecting the weights to minimize $D_{CM}(p)$ naturally produces density estimates that are closer to $f(x)$ than are estimates from minimizing $D_{L_2}(p)$.

Note that using the power-divergence metric (Cressie and Read (1984)),

$$D_\rho(p) = \frac{1}{\rho(1-\rho)} \left(n - \sum_{i=1}^n (np_i)^\rho \right),$$

in this setting may not be useful, because it requires that the p_i used to estimate a density from (2.2) be nonnegative (p_i must satisfy $\sum_{i=1}^n p_i = 1$ and $p_i \geq 0$ for this approach), and some constraints may require negative weights. Furthermore, although $D_\rho(p)$ has an appealing immediate interpretation as a measure of entropy, it does require that the user select an additional tuning parameter for its implementation (ρ). Lastly, as Hall and Huang (2002) note, problems can arise as p_i approaches zero, because enforcing constraints on a curve leads to “data compression” (i.e., the effective sample size used locally is smaller than the corresponding effective sample size for the unconstrained estimator). This difference is achieved by setting some of the constraint weights to zero. This information is not lost however, but simply reassigned to observations that receive nonzero weights. Thus, there can be substantial differences between our elected metrics and $D_\rho(p)$; while both $D_{L_2}(p)$ and $D_{CM}(p)$ behave well when p_i approaches zero, $D_\rho(p)$ may not be applicable for certain constraints with particular values of ρ .

2.2. Bounded support PDF kernel functions

We wish to develop an approach that will suit the many and varied needs of a range of practitioners. *Boundary bias* affects the quality of kernel density estimates when substantial probability mass occurs at a support boundary. The most well-known solutions to this problem are *data-reflection*, *data-transformation*, and *kernel carpentry*. Data-reflection involves duplicating data symmetrically (i.e., reflecting) around its boundary, running standard bandwidth selection and kernel estimation, and then adjusting the resulting estimate to ensure it is proper (i.e., integrates to one) on its support. Data-transformation involves some mathematical transform of the data that, when rescaled, has the desired effect. Kernel carpentry uses kernel functions that adapt to the presence of a boundary, thus mitigating the effect of the boundary. To

some degree, these methods all reduce the amount of bias that would otherwise be present near a boundary to that which holds in the interior of the support, where it is free from boundary effects (in effect, lying h or greater distance from the boundary in the interior). However, data-reflection and transformation require extra steps of the user, which is both inconvenient and unnecessary. In what follows, we take a kernel carpentry approach, and adopt truncated kernel functions of the type

$$K(z, a, b) = \begin{cases} \frac{K(z)}{G(z_b) - G(z_a)} & \text{if } z_a \leq z \leq z_b, \\ 0 & \text{otherwise,} \end{cases}$$

where $z = (x - X)/h$, with X the random variable representing X_i , $z_b = (b - x)/h$, $z_a = (a - x)/h$, and $G(z) = \int_{-\infty}^z K(t) dt$. Given that $K(z)$ is a standard univariate kernel function, $G(z)$ is the CDF counterpart to the PDF $K(z)$ that we used to estimate $F(x)$. Note that if $K(z)$ is, for instance, the Gaussian density function, then $K(z, a, b)$ is simply the (doubly) truncated Gaussian density function. When $a = -\infty$ and $b = \infty$, then $K(z, a, b) = K(z)$, which is a standard kernel function, such as the Gaussian (or Epanechnikov). Hence this kernel function allows for unbounded or compact support without modification. When conducting a constrained estimation, it may be necessary to use the integrated version of $K(z, a, b)$, or derivatives thereof. We briefly outline some helpful relationships used to obtain these objects from the doubly truncated kernel function $K(z, a, b)$.

2.2.1. Integral kernel functions (e.g., CDF kernels)

To reduce the notational burden, let $H_{ba}(z) = H(z_b) - H(z_a)$, for any function $H(\cdot)$. To estimate a CDF using kernel methods in the presence of support bounds, we can obtain the counterpart to $K(z, a, b)$ by adopting the following transformation for (doubly) truncated density functions:

$$G(z, a, b) = \begin{cases} 0 & \text{if } z < z_a, \\ \frac{G(\max(\min(z, z_b), z_a)) - G(z_a)}{G_{ba}(z)} & \text{if } z_a \leq z \leq z_b, \\ 1 & \text{otherwise.} \end{cases}$$

2.2.2. Derivative kernel functions

Some of the constraints we consider are placed on the derivative of the kernel density estimates, and hence we may require derivatives of the kernel function. To that end, we apply the quotient rule to obtain the first derivative of the doubly truncated kernel function, yielding

$$K'(z, a, b) = \begin{cases} \frac{K'(z)}{G_{ba}(z)} - \frac{K(z)G'_{ba}(z)}{G_{ba}(z)^2} & \text{if } z_a \leq z \leq z_b, \\ 0 & \text{otherwise.} \end{cases}$$

Note that

$$\frac{K(z)G'_{ba}(z)}{G_{ba}(z)^2} = K(z, a, b) \frac{K_{ba}(z)}{G_{ba}(z)}.$$

The second derivative is found by applying the quotient and the product rules, yielding

$$K''(z, a, b) = \begin{cases} \frac{d}{dx} \frac{K'(z)}{G_{ba}(z)} - \frac{d}{dx} \frac{K(z)G'_{ba}(z)}{G_{ba}(z)^2} & \text{if } z_a \leq z \leq z_b, \\ 0 & \text{otherwise.} \end{cases}$$

Note that the first term on the right-hand side can be expressed as

$$\frac{d}{dx} \frac{K'(z)}{G_{ba}(z)} = \frac{K''(z)}{G_{ba}(z)} - \frac{K'(z)K_{ba}(z)}{G_{ba}(z)^2},$$

and the second term (ignoring the minus sign) can be expressed as

$$\begin{aligned} \frac{d}{dx} \frac{K(z)K_{ba}(z)}{G_{ba}(z)^2} &= \frac{K'(z)(K(z_b) - K(z_a)) + K(z)K'_{ba}(z)}{G_{ba}(z)^2} \\ &\quad - \frac{2K(z)K_{ba}(z)G_{ba}(z)G'_{ba}(z)}{G_{ba}(z)^4} \\ &= \frac{K'(z)K_{ba}(z) + K(z)K'_{ba}(z)}{G_{ba}(z)^2} - \frac{2K(z)K_{ba}(z)^2}{G_{ba}(z)^3}. \end{aligned}$$

Therefore, we obtain

$$K''(z, a, b) = \begin{cases} \frac{K''(z)}{G_{ba}(z)} - \frac{2K'(z)K_{ba}(z) + K(z)K'_{ba}(z)}{G_{ba}(z)^2} + \frac{2K(z)K_{ba}(z)^2}{G_{ba}(z)^3} & \text{if } z_a \leq z \leq z_b, \\ 0 & \text{otherwise.} \end{cases}$$

Note that, for the Gaussian kernel, if $a = -\infty$ and $b = \infty$, then $K(z_a) = K(z_b) = K'(z_a) = K'(z_b) = 0$, and $G(z_b) - G(z_a) = 1$; hence, $K'(z, a, b) = K'(z)$ and $K''(z, a, b) = K''(z)$ in the unbounded support case, as expected.

The utility of this doubly truncated kernel function is that it can directly admit unbounded support (i.e., on $(-\infty, \infty)$), support on $[a, \infty)$ with a finite, support on $(-\infty, b]$ with b finite, and support on $[a, b]$ with both a and b finite, without further modification. Using this kernel function allows us to deliver an approach that directly admits support bounds *and* shape constraints, which we believe enhances its practical appeal by increasing its potential application.

2.3. Hypothesis testing

We can test the validity of the shape constraints being imposed by following the insights of Hall et al. (2001) and Du, Parmeter and Racine (2013), and using a bootstrap inferential procedure. Briefly, the test statistic is the value of the objective function from solving the quadratic program when imposing the constraints. The bootstrap procedure draws bootstrap resamples from the

null (i.e., constrained) density in order to construct the null distribution of the test statistic (i.e., the value of the objective function from solving the quadratic program when imposing the constraints on the bootstrap resamples). The test involves computing a P -value constructed by comparing the test statistic with that obtained from the empirical distribution constructed under the null or, alternatively, by comparing the test statistic with the desired $1 - \alpha$ quantile obtained from the empirical null distribution where α is the desired size of the test procedure (the test is one-sided with a right-tailed rejection region).

More specifically, this bootstrap approach involves estimating the constrained density $\hat{f}(\mathbf{x}|p)$ based on the sample realizations $\{\mathbf{X}_i\}$; and then rejecting H_0 if the observed value of $D_j(\hat{p})$ is too large, where $j \in \{L_2, CM\}$. To ensure that the constraints are satisfied, we propose sampling from $\hat{f}(\mathbf{x}|p)$ rather than from $\hat{f}(\mathbf{x}|p_{unif})$. A simple way to do this is to use rejection sampling.

These resamples are generated under H_0 . Hence we recompute $\hat{f}(\mathbf{x}|p)$ for the bootstrap sample $\{\mathbf{X}_i^*\}$; which we denote as $\hat{f}(\mathbf{x}|p^*)$, yielding $D_j(p^*)$. We repeat this process B times. Finally, we compute the empirical P value, P_B , which is simply the proportion of the B bootstrap resamples $D_j(p^*)$ that exceed $D_j(\hat{p})$, that is,

$$P_B = 1 - \hat{F}(D_j(\hat{p})) = \frac{1}{B} \sum_{j=1}^B I(D_j(p^*) > D_j(\hat{p})),$$

where $I(\cdot)$ is the indicator function and $\hat{F}(D_j(\hat{p}))$ is the empirical distribution function of the bootstrap statistics. Then, we reject the null hypothesis if P_B is less than α , the level of the test.

We now consider a few illustrative applications of imposing shape restrictions, before turning to the theoretical underpinnings of the proposed method.

2.4. Illustrative applications: Monotonicity and concavity

Monotonicity and concavity constraints are two popular shape constraint domains that our approach can cover. As in Du, Parmeter and Racine (2013), we solve a simple quadratic program using (2.6) to generate the constrained estimate. Figure 1 presents the results for a bounded density on $[0, 1]$ imposing monotonicity (the distribution is Beta(5,1)). For this simple illustration, we generate 100 observations and select the bandwidth using Silverman's rule-of-thumb approach. We see little difference between the constrained and the unconstrained estimators for $x > 0.6$; all of the constraint enforcement occurs in the left tail of the density. Given our restriction that the weights sum to zero, this leads to only minor changes in the shape of the density beyond where the constraints need to be enforced. This becomes clearer by looking at the lower plot in Figure 1, which plots the constrained and unconstrained derivative estimates.

Figure 2 presents the results when imposing concavity on an unbounded support random variable (the distribution is $N(0, 1)$). Once again, we generate

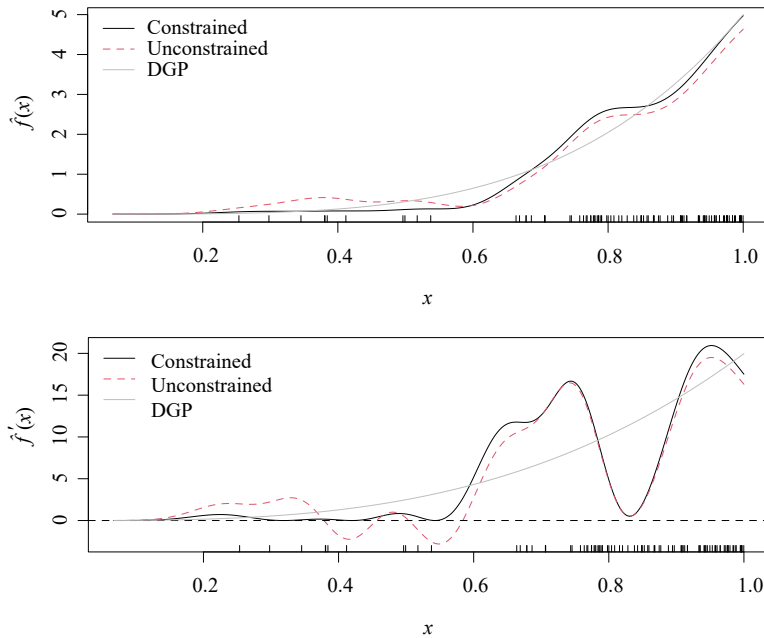


Figure 1. Monotone shape-constrained density estimation ($\hat{f}'(x) \geq 0$). The upper figure plots the constrained and unconstrained density estimates, the lower figure plots the constrained and unconstrained first derivative estimates.

100 observations randomly and construct the bandwidth using Silverman's rule-of-thumb. Here, we enforce concavity on the density, which is *not* a property of the Gaussian density (though it *is* log-concave). We see that enforcing *invalid* constraints produces substantial distortions in both the density and the corresponding first derivative, as expected.

2.5. Log-concave kernel density estimation

Log-concavity is a popular constraint, although it is only one of many shape constraint domains that our approach can cover. To impose log-concavity/convexity, we require $d^2 \log(\hat{f}(x))/dx^2$ and $d^2 K(Z_i)/dx^2$. The former is given by

$$\frac{d^2 \log(\hat{f}(x))}{dx^2} = \frac{\hat{f}''(x)\hat{f}(x) - (\hat{f}'(x))^2}{(\hat{f}(x))^2},$$

and the latter is given by

$$\frac{d^2 K(Z_i)}{dx^2} = K''(Z_i).$$

Note that

$$\hat{f}'(x) = \frac{1}{nh} \sum_{i=1}^n K'(Z_i),$$

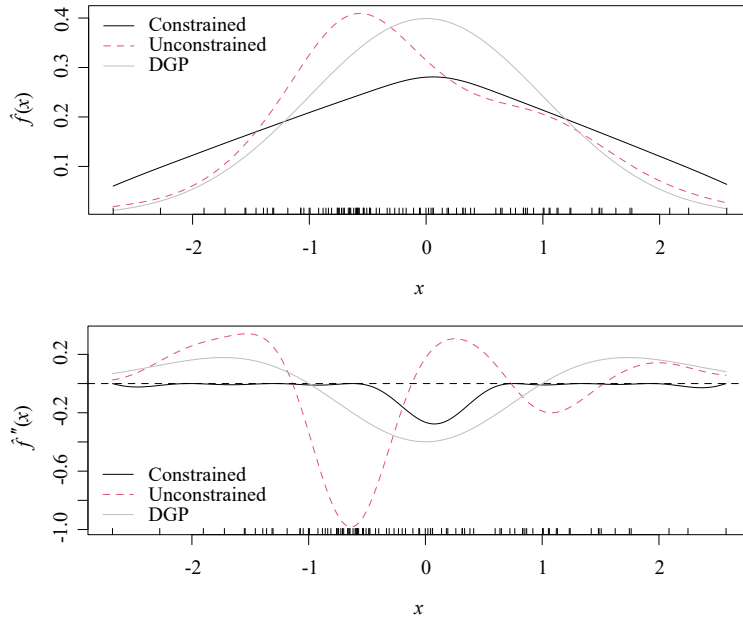


Figure 2. Concave shape-constrained density estimation ($\hat{f}''(x) \leq 0$). The upper figure plots the constrained and unconstrained density estimates, the lower figure plots the constrained and unconstrained second derivative estimates.

$$\hat{f}''(x) = \frac{1}{nh} \sum_{i=1}^n K''(Z_i).$$

2.6. Illustrative application: Log-concavity

Figure 3 presents the results for a draw from the $N(0, 1)$ Gaussian distribution. The Gaussian density is log-concave, but the kernel estimate need not be, as the following example illustrates. We generate 250 observations from a standard normal distribution, and use Silverman's rule-of-thumb bandwidth to smooth the density. As in Figure 1, there is little difference between the constrained and the unconstrained estimates. Moreover, the log-densities are also quite similar, aside from one region of nonconcavity of the log-density for $-2.5 < x < -1.9$. Both the constrained and the unconstrained densities integrate to one and are proper.

2.7. Categorical (ordered) PMFs

The approach we consider for a shape-constrained PDF estimation can also be applied to a shape-constrained PMF estimation (Aitchison and Aitken (1976); Racine, Li and Yan (2020)). When X is an ordered categorical variable ($X \in \mathbb{D} = \{D_0, D_1, \dots, D_{c-1}\}$, where c is the number of (ordered) outcomes), we need only the one value of a_i per outcome (because $a_i = a_j$ when $X_i = X_j$). When placing shape constraints on derivatives, we adopt the classical convention that

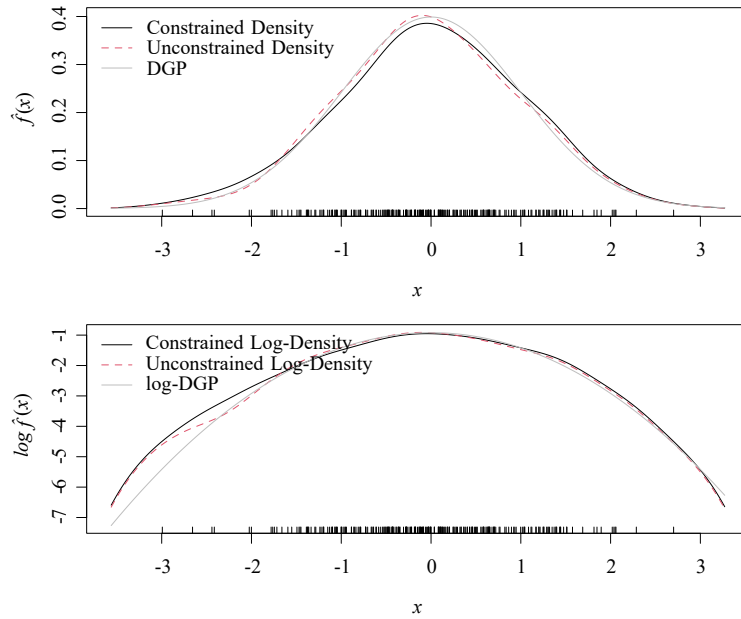


Figure 3. Log-concave shape-constrained density estimation. The upper figure plots the constrained and unconstrained density estimates, the lower figure plots the constrained and unconstrained log-density estimates.

for discrete support variables, derivatives are defined in terms of simple finite differences. For an ordered discrete random variable, we use the notation $P(x) = Pr(X = x)$ to denote the PMF. Let $\hat{P}(x)$ denote the kernel estimate of $P(x)$ given by

$$\hat{P}(x) = \frac{1}{n} \sum_{i=1}^n l(X_i, x, \lambda),$$

where $l(X_i, x, \lambda)$ is an appropriate kernel function for ordered discrete support random variables. The counterpart of the first derivative in this setting is $\Delta_j(x) = (P(x_{(j)}) - P(x_{(j-1)})) / (x_{(j)} - x_{(j-1)})$, where $x_{(j)}$ are the order statistics, that can be computed directly from an unconstrained estimate (as can higher-order derivatives, if needed). As was the case for the shape-constrained PDF estimation, the counterpart to (2.3) for the PMF estimation can be written as

$$\hat{P}(x|p) = \hat{P}(x) + \sum_{i=1}^n a_i l(X_i, x, \lambda), \quad (2.7)$$

where λ is the smoothing parameter analogous to the bandwidth h for its continuous support counterpart. The mechanics of the shape-constrained PMF estimator are the same as those for the shape-constrained PDF estimation described previously, and so are not repeated here (see Racine, Li and Yan (2020) for further details). We now consider an empirical illustration based on count data

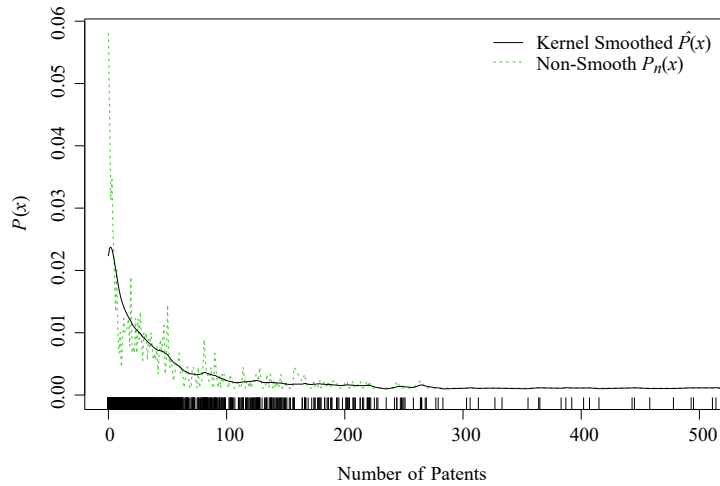


Figure 4. Unconstrained smooth and nonsmooth PMF estimates for patent data. The smooth estimate appears as a solid line, the nonsmooth estimate as a dotted line.

that have ordered discrete support.

2.8. Empirical application: Shape-constrained PMFs

We consider a data set collected by Hausman, Hall and Griliches (2002) that records the number (count) of successful patent applications by 128 U.S. firms across a seven-year period (1968–1974). We model the kernel-smoothed PMF for the number of successful patent applications with likelihood cross-validated bandwidth selection, and present the results in Figure 4. The nonsmooth estimate is quite noisy, whereas the smooth estimate is much less so. Like its empirical counterpart, the smooth estimate delivers probability estimates that *sum* to one, but the smooth estimate is expected to be more efficient from a squared error perspective.

Figure 4 reveals that the unconstrained kernel PMF estimator, though perhaps more plausible an estimate than the nonsmooth empirical estimator, implausibly changes sign in many places. A perhaps more reasonable assumption is that the estimate is monotonically decreasing. Hence we consider imposing this shape constraint on the kernel PMF estimate. Figure 5 presents the smooth unconstrained and monotonically constrained estimates. As noted above, the derivatives for the PMF estimate are given by $\Delta_j(x) = (P(x_{(j)}) - P(x_{(j-1)})) / (x_{(j)} - x_{(j-1)})$, where $x_{(j)}$ are the order statistics, which can be computed directly. The weight matrix required to solve the quadratic program is then the difference between kernel functions evaluated at $x_{(j)}$ and $x_{(j-1)}$ divided by the difference $x_{(j)} - x_{(j-1)}$. To impose the monotonically decreasing constraint, we define $\Delta_1(x) \leq 0$ (we reverse this definition for monotonically increasing constraints).

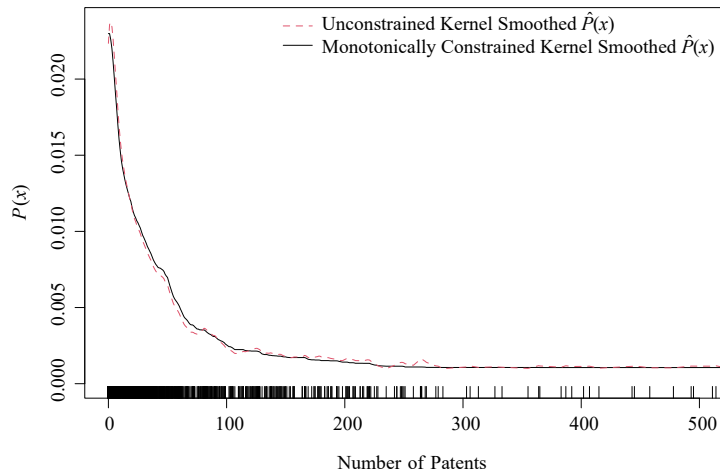


Figure 5. Unconstrained and constrained smooth probability function estimates for the patent data.

3. Theoretical Properties of the Constrained Estimator

In this section, we provide four key theoretical results. First, under weak conditions, the constraint weights generated by our approach are shown to be well defined and unique. Second, we demonstrate the consistency of the constrained density estimator, where appropriate, in terms of its closeness to the unconstrained density estimator, which is well known to be consistent. We consider three distinct settings: (i) when the constraints are indeed true on the entire support of X ; (ii) when the constraints are satisfied everywhere except at points of measure zero; and (iii) when the constraints are violated on a set with positive measure. For (i) and (ii), we establish the consistency of the constrained density estimator under weak conditions on the order of the derivatives of the true density and on the bandwidth (naturally, (iii) does not allow for consistent estimation). Third, we extend our results in the continuous case to those for ordered PMFs. Here, we are only able to establish consistency when the constraints hold on the entire support of the discrete random variable; nevertheless, these results are novel and of practical value. Fourth, we provide the asymptotic distribution of our proposed test statistic when testing the null hypothesis of the validity of the shape constraints being imposed.

Our theoretical results for continuous data are similar to those of Hall et al. (2001) and Du, Parmeter and Racine (2013), but with four important differences. First, Hall et al. (2001) and Du, Parmeter and Racine (2013) impose constraints in a regression setting. The density setting is complicated by the lack of an error term, such that we cannot apply existing theory directly. Second, Hall et al. (2001) use the power-divergence measure of Cressie and Read (1984) and Du, Parmeter and Racine (2013) use the L_2 -metric. Here, we establish the

consistency of the constrained estimator using the objective function proposed by Li, Liu and Li (2017), which, rather than selecting constraint weights as close as possible to the uniform weights (as in Du, Parmeter and Racine (2013)), selects weights as close as possible to the unconstrained estimator. Intuitively, this modification makes sense, given that the unconstrained estimator is consistent to begin with. Although Li, Liu and Li (2017) show the impressive finite-sample properties of their objective function when selecting constraint weights for a constrained K_{nn} regression estimator, the change in the objective function also necessitates changes to existing theory. Third, existing theory works relatively well for constraints on the density. However, several additional modifications are required when imposing constraints on, for example, the log-density. Fourth, we develop the appropriate theory for the constrained estimation of the PMF. To the best of our knowledge, this is the first application of these types of constrained methods to kernel-smoothed discrete data.

To begin, \mathbf{X}_i is of dimension r . Our goal is to impose constraints on the density (or log-density) of the form $f^{(\mathbf{s})}(\mathbf{x}) = [\partial^{s_1} f(\mathbf{x}) \cdots \partial^{s_r} f(\mathbf{x})] / [\partial x_1^{s_1} \cdots \partial x_r^{s_r}]$ (or $\log f^{(\mathbf{s})}(\mathbf{x})$), where \mathbf{s} is an r -vector corresponding to the dimension of \mathbf{x} . Note that the general two-sided constraints in (2.5) can be expressed as one-sided constraints of the form

$$\sum_{\mathbf{s} \in \mathbf{S}_k} \alpha_{\mathbf{s},k} f^{(\mathbf{s})}(\mathbf{x}) - c_k(\mathbf{x}) \geq 0, \quad k = 1, \dots, T, \tag{3.1}$$

where T is the total number of restrictions, with the sum taken over all density derivative vectors in \mathbf{S}_k , and $\alpha_{\mathbf{s},k}$ is used to generate the appropriate constraints imposed on the density derivatives ($j = 0, 1, \dots$). This notation admits an arbitrary number of internally consistent constraints imposed simultaneously on the density and its derivatives, though in most cases, we expect that a single constraint (i.e., $T = 1$) will suffice. As an example, for $r = 1$ and the imposition of monotonicity, we have $T = 1$ with $\mathbf{s} = (1)$, $\mathbf{S}_k = \{(1)\}$, $\alpha_{\mathbf{s},k} = 1$, and $c_k(\mathbf{x}) = 0$, for all \mathbf{x} .

Before formally developing the theory for our general constrained density estimator, we introduce some additional simplifying notation. Denote the domain of interest by $\mathcal{J} \equiv [\mathbf{m}, \mathbf{b}] = \prod_{i=1}^r [m_i, b_i]$. We also define a differential operator $f \mapsto f^{\mathcal{D}}$ such that $f^{\mathcal{D}}(\mathbf{x})$ is a length- T , vector with k th entry $\sum_{\mathbf{s} \in \mathbf{S}_k} \alpha_{\mathbf{s},k} f^{(\mathbf{s})}(\mathbf{x})$. We take $|\mathbf{s}| = \sum_{i=1}^r s_i$ as the *order* for a derivative vector $\mathbf{s} = (s_1, \dots, s_r)$, and say a derivative \mathbf{s}_1 has a *higher order* than that of \mathbf{s}_2 if $|\mathbf{s}_1| > |\mathbf{s}_2|$. Let $\mathbf{S} = \cup_{k=1}^T \mathbf{S}_k$ and $\mathbf{d}_{\mathbf{S}}$ be the derivative of the *maximum order* among all the derivatives in \mathbf{S} ; for simplicity, we drop the subscript \mathbf{S} from $\mathbf{d}_{\mathbf{S}}$. Without loss of generality, we set $c_k(x) = 0$ in what follows. Plugging (2.2) into (3.1) yields

$$\sum_{i=1}^n p_i K_i^{\mathcal{D}}(\mathbf{x}) \geq 0. \tag{3.2}$$

Here, $K_i^{\mathcal{D}}(\mathbf{x})$ represents the form of the constraints based on the appropriate kernel derivatives, that is, it subsumes the appropriate entries of the derivative vector $f^{\mathcal{D}}(\mathbf{x})$. Lastly, we define $\tilde{f}(x) = \hat{f}(x|p_{unif})$ to further simplify our notation.

Although the theory we present here is capable of imposing constraints on either the density or the log-density, for notational simplicity, we presume that the practitioner is interested in only one or the other.

3.1. Existence of the constrained PDF estimator

The first result that we establish is an existence result, that is, that a set of weights exists, provided that the constraints imposed are internally consistent and satisfy the constraints in (3.2).

Theorem 1 (Existence). *Assume that the set $\{1, \dots, n\}$ contains a sequence $\{i_1, \dots, i_k\}$ with the following properties:*

- i) for each $\ell = 1, \dots, k$, $K_{i_\ell}^{\mathcal{D}}(\mathbf{x})$ is strictly positive and continuous on an open set $\mathbf{O}_{i_\ell} \subset \mathbb{R}^r$, and vanishes on $\mathbb{R}^r \setminus \mathbf{O}_{i_\ell}$;*
- ii) every $\mathbf{x} \in \mathcal{J}$ is contained in at least one open set \mathbf{O}_{i_k} ;*
- iii) for $1 \leq \ell \leq n$, $K_{i_\ell}^{\mathcal{D}}(\mathbf{x})$ is continuous on $(-\infty, \infty)^r$.*

Then, there exists a vector $p = (p_1, \dots, p_n)$ such that the constraints are satisfied for all $\mathbf{x} \in \mathcal{J}$.

Conditions *i)* and *ii)* of Theorem 1 ensure the existence of an open cover of the domain \mathcal{J} by the open sets \mathbf{O}_{i_ℓ} on which $K_{i_\ell}^{\mathcal{D}}$ is positively supported for some i_ℓ . Note that the above conditions are sufficient, but not necessary for the existence of a set of weights that satisfy the constraints for all $\mathbf{x} \in \mathcal{J}$. For example, if $\text{sign } K_{j_n}^{\mathcal{D}}(\mathbf{x}) = 1 \ \forall \mathbf{x} \in \mathcal{J}$ for some sequence j_n in $\{1, \dots, n\}$, and $\text{sign } K_{l_n}^{\mathcal{D}}(\mathbf{x}) = -1 \ \forall \mathbf{x} \in \mathcal{J}$ for another sequence l_n in $\{1, \dots, n\}$, then for those observations that switch signs, p_i may be set equal to zero, and $p_{j_n} > 0$ and $p_{l_n} < 0$ are sufficient to ensure the existence of a set of p satisfying the constraints. The proof of Theorem 1 is provided in the Supplemental Material.

3.2. Consistency of the constrained PDF estimator

Here, we discuss the consistency of our constrained estimator. To begin, define a *hyperplane subset* of \mathcal{J} as a subset of the form $\mathcal{S} = \{x_{0k} \times \prod_{i \neq k} [m_i, b_i]\}$, for some $1 \leq k \leq r$ and some $x_{0k} \in [m_k, b_k]$. We call \mathcal{S} an *interior hyperplane subset* if $x_{0k} \in (m_k, b_k)$. In the following, $f(\cdot)$ (or $f^{\mathcal{D}}(\cdot)$) is the true density (or its derivative), \hat{p} is the optimal weight vector satisfying the constraints, $\hat{f}(\cdot|\hat{p})$ (or $\hat{f}^{\mathcal{D}}(\cdot|\hat{p})$) is the constrained estimator defined in (2.3), and $\tilde{f}(\cdot)$ (or $\tilde{f}^{\mathcal{D}}(\cdot)$) is the unconstrained estimator defined in (2.3).

Assumption A1.

- i) The sample \mathbf{X}_i either forms a regularly spaced grid on a compact set $\mathcal{I} \equiv [\mathbf{c}, \mathbf{e}] = \prod_{i=1}^r [c_i, e_i]$, or constitutes independent random draws from a distribution with a density f that is continuous and nonvanishing on \mathcal{I} ; the kernel function $K(\cdot)$ is a symmetric, compactly supported density such that $K^{\mathcal{D}}$ is Hölder-continuous on $\mathcal{J} \subset \mathcal{I}$.
- ii) $f^{\mathcal{D}}$ is continuous on \mathcal{J} .
- iii) The bandwidth associated with each variable, h_j , satisfies $h_j \propto n^{-1/(3r+2|\mathbf{d}|)}$, for $1 \leq j \leq r$, where $|\mathbf{d}|$ is the maximum order of the derivative vector \mathbf{d} .
- iv) The true density f is bounded away from zero, say, $f(\mathbf{x}) > \tau$, for some fixed constant $\tau > 0$.

Assumption A1 i) is standard in the kernel density literature; at the expense of a more tedious proof, the same results can be demonstrated if the density is assumed to exist on an r -dimensional ball instead of on a hypercube. Assumption A1 ii) ensures the requisite smoothness of $f^{\mathcal{D}}$. Note that the bandwidth rate in Assumption A1 iii) is, in general, higher than the standard optimal rate $n^{-1/(r+4)}$. However, this is not surprising for our restricted problem. The optimal rate only guarantees the convergence of our unrestricted function estimator \tilde{f} . However, the restricted problem also requires the convergence of the derivative $\tilde{f}^{\mathcal{D}}$, which often needs a higher bandwidth rate. In the single-predictor monotone regression problem considered in Hall et al. (2001), this rate happens to coincide with the optimal rate $n^{-1/5}$. Furthermore, when the bandwidths all share the same rate, one can rescale each component of \mathbf{x} to ensure a uniform bandwidth $h \propto n^{-1/(3r+2|\mathbf{d}|)}$ for all components. This simplification is made without loss of generality. Thus, we use h^r rather than $\prod_{j=1}^r h_j$, for notational simplicity. If we consider densities on a compact interval, then Assumption A1 iv) is not so restrictive. However, it may not work for common densities, such as the normal and exponential densities.

Theorem 2 (Consistency). *Suppose that Assumption A1 1.–4. holds.*

- i) *If $f^{\mathcal{D}} > 0$ on \mathcal{J} , then, with probability one, $\hat{p} = 1/n$ for all sufficiently large n , and $\hat{f}^{\mathcal{D}}(\cdot|\hat{p}) = \tilde{f}^{\mathcal{D}}$ on \mathcal{J} for all sufficiently large n . Hence, $\hat{f}(\cdot|\hat{p}) = \tilde{f}$ on \mathcal{J} for all sufficiently large n .*
- ii) *Suppose that $f^{\mathcal{D}} > 0$, except on an interior hyperplane subset $\mathcal{X}_0 \subset \mathcal{J}$, where we have $f^{\mathcal{D}}(\mathbf{x}_0) = 0, \forall \mathbf{x}_0 \in \mathcal{X}_0$. In addition, for any $\mathbf{x}_0 \in \mathcal{X}_0$, suppose that $f^{\mathcal{D}}$ has second-order continuous derivatives in the neighborhood of \mathbf{x}_0 , with $\partial f^{\mathcal{D}}/\partial \mathbf{x}(\mathbf{x}_0) = \mathbf{0}$ and $\partial^2 f^{\mathcal{D}}/\partial \mathbf{x} \partial \mathbf{x}^T(\mathbf{x}_0)$ nonsingular; then, $|\hat{f}(\cdot|\hat{p}) - \tilde{f}| = O_p(h^{|\mathbf{d}|+(r+1)/2})$ uniformly on \mathcal{J} .*

iii) Under the conditions in ii), there exist random variables $\Theta = \Theta(n)$ and $Z_1 = Z_1(n) \geq 0$ satisfying $\Theta = O_p(h^{|\mathbf{d}|+r+1})$ and $Z_1 = O_p(1)$, such that $1 - \Theta \leq \hat{f}(x|\hat{p})/\tilde{f}(x) \leq 1 + \Theta$ uniformly for $\mathbf{x} \in \mathcal{J}$, with $\inf_{\mathbf{x}_0 \in \mathcal{X}_0} |\mathbf{x} - \mathbf{x}_0| > Z_1 h^{(r+1)/4}$.

In Theorem 2, part i) suggests that when the constraint is strictly satisfied by the true function, the constrained estimator $\hat{f}(\cdot|\hat{p})$ and the unconstrained estimator \tilde{f} are essentially the same, and thus share the same rate of convergence. Part ii) gives the order of difference between $\hat{f}(\cdot|\hat{p})$ and \tilde{f} when $f^{\mathcal{D}} = 0$ on an interior hyperplane. Note that the order in ii) indicates a different convergence rate of $\hat{f}(\cdot|\hat{p})$ from that of \tilde{f} in such a case. Part iii) is concerned with the asymptotic behavior of the weights \hat{p} in such a case. Note that these results are easily extendable to the case of $f^{\mathcal{D}} \leq 0$ with a switch of sign in f .

The proof of Theorem 2 appears in the online Supplementary Material. Theorem 2 is a multivariate, multi-constraint, hyperplane subset adaptation of Du, Parmeter and Racine (2013) to density estimation using the metric in (2.6).

3.3. Theoretical properties of the constrained PMF estimator

Theorems 1 and 2 can be extended to the ordered discrete support setting under similar assumptions, though with some important modifications required.

Assumption B1.

- i) Assume that the set $\{1, \dots, n\}$ contains a sequence $\{i_1, \dots, i_k\}$ with the following properties:
 - (i) for each k , $\ell_{i_k}^{\mathcal{D}}(x)$ is strictly positive on a nonempty set $\mathbf{O}_{i_k} \subset \mathbb{D}$, and vanishes on $\mathbb{D} \setminus \mathbf{O}_{i_k}$; (ii) every $x \in \mathbb{D}$ is contained in at least one nonempty set \mathbf{O}_{i_k} .
- ii) Assume that the kernel function $l(\cdot)$ in (2.7) is an ordered kernel function, and that the smoothing parameter λ in (2.7) is of order $\lambda = O_p(n^{-1})$, which is a standard result in the literature.

Assumption B1 i) is similar to the sufficient conditions in Theorem 1 for the continuous case. For the smoothing parameter condition in Assumption B1 ii), Ouyang, Li and Racine (2006) show that a smoothing parameter λ selected using cross-validation can have order $O_p(n^{-1})$, as long as the marginal distributions of X are not all uniform.

Theorem 3 (PMF Estimator). *Suppose that Assumption B1 holds. Then, we have the following properties for the constrained PMF estimate $\hat{P}^{\mathcal{D}}(\cdot|\hat{p})$. Our use here of the differential is with respect to the difference order, as opposed to differentiation.*

- i) There exists a vector $p = (p_1, \dots, p_n)$ such that the constraints are satisfied for all $x \in \mathbb{D}$.
- ii) If $P^{\mathcal{D}} > 0$ on \mathbb{D} then, with probability one, $\hat{p} = 1/n$ for all sufficiently large n , and $\hat{P}^{\mathcal{D}}(\cdot|\hat{p}) = \tilde{P}^{\mathcal{D}}$ on \mathbb{D} for all sufficiently large n . Hence, $\hat{P}(\cdot|\hat{p}) = \tilde{P}$ on \mathbb{D} for all sufficiently large n .

The proof of existence requires only minor changes to the proof for the continuous data setting, and is thus omitted. The proof for consistency still requires taking differences across the cells of the discrete random variable, which suggests that our constraints correspond to an ordered discrete random variable (Li and Racine (2002)).

For the proof of the consistency, note that parts ii) and iii) cannot be generalized. This result has a straightforward intuitive explanation. In the continuous-only setting, these parts focus on the case where the constraint is violated on a set of measure zero. The argument is that, even if the constraint is violated, as long as it occurs on an interior subset hyperplane, the constrained estimator is still a consistent estimator for the unknown density (except on a set of measure zero). In the discrete data setting, these results no longer hold, because for a discretely supported random variable, a measure-zero event is equivalent to an outcome not in the support; thus, a violation of the constraint is more troubling when considering discrete data.

3.4. Asymptotic properties of $D(\hat{p})$ asymptotic properties of $D(\hat{p})$

Our discussion on inference of the smoothness constraints follows the same setup as in Du, Parmeter and Racine (2013). We focus on using the L_2 -norm rather than CM, because a closed-form solution for the optimal weights is mathematically more tedious, owing to the cross-products of the weights in the objective function. Note that the asymptotic expansions of the weights between L_2 and CM are of the same order, but will obviously be of a slightly different form. Let $\psi_i(\mathbf{x}) = K_i^{\mathcal{D}}(\mathbf{x})$, for $i = 1, \dots, n$.

Recall that our minimization problem is

$$\min_{p_1, \dots, p_n} \sum_{i=1}^n (n^{-1} - p_i)^2, \quad \text{s.t.} \quad \sum_{i=1}^n p_i = 1, \quad \sum_{i=1}^n p_i \psi_i(\mathbf{x}) \geq 0, \forall \mathbf{x}.$$

In practice, this minimization is carried out by taking a fine grid $(\mathbf{x}_1, \dots, \mathbf{x}_N)$, where N is large, and solving

$$\min_{p_1, \dots, p_n} \sum_{i=1}^n (n^{-1} - p_i)^2, \quad \text{s.t.} \quad \sum_{i=1}^n p_i = 1, \quad \sum_{i=1}^n p_i \psi_i(\mathbf{x}_j) \geq 0, 1 \leq j \leq N. \quad (3.3)$$

We place the same assumption on the grid points $(\mathbf{x}_1, \dots, \mathbf{x}_N)$ as in Du, Parmeter and Racine (2013).

Assumption B2.

i) $N \rightarrow \infty$ as $n \rightarrow \infty$ and $N = O(n)$.

ii) Let $d_N = \inf_{1 \leq j_1, j_2 \leq N} |\mathbf{x}_{j_1} - \mathbf{x}_{j_2}|$ be the minimum distance between grid points. We require $d_N \rightarrow 0$ and $h^{-1}d_N \rightarrow \infty$.

Assumption B2 essentially dictates how the grid points behave. We need to ensure that the grid becomes effectively dense as n increases (Assumption B2 i)), while also needing the speed at which the smallest distance between the grid points decays to be slower than the rate of decay of the smoothing parameters (Assumption B2 ii)). The latter assumption is necessary to eliminate correlation across $\psi_i(\mathbf{x})$ as n grows (Chacón, Duong and Wand (2011)).

Let \hat{p}_i , for $i = 1, \dots, n$, be the solution to the quadratic programming problem in (3.3). Then, the asymptotic distribution of $D(\hat{p})$ is given in the following theorem, with the proof given in the Supplementary Material.

Theorem 4. *Suppose that assumptions A1 i)–iv) and B1 i)–iv) hold. Then, as $n \rightarrow \infty$, we have*

$$\frac{n^2 \sigma_{K^{(a)}}^2}{h^{2|d|+r} \left(\sum_{j=1}^M f^{\mathcal{D}}(\mathbf{x}_j^*) \right)^2} D(\hat{p}) \sim \chi^2(n), \quad (3.4)$$

where $\sigma_{K^{(a)}}^2 = \int [K^{(d)}(y)]^2 dy$, and $\{\mathbf{x}_1^*, \dots, \mathbf{x}_M^*\} \subset \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ are the slack points defined in the Supplementary Material.

Theorem 4 is the density equivalent of the regression-based test proposed by Du, Parmeter and Racine (2013). Aside from several structural details, the main result follows from their initial theory. The diverging degrees of freedom is expected, because H_0 and H_1 are both evaluated on infinite-dimensional parameter spaces (see also Fan, Zhang and Zhang (2001)). Note too that, similarly to the generalized likelihood ratio test statistic of Fan, Zhang and Zhang (2001), the distributional convergence in (3.4) is equivalent to $\sqrt{2n}(T_n - n) \xrightarrow{\mathcal{L}} N(0, 1)$, where T_n is the statistic on the left-hand side of (3.4).

Given the well-known issues with the speed of convergence of nonparametric tests, we recommend using a bootstrap algorithm instead. Another reason to prefer the bootstrap is that the normalizing constant in (3.4) requires that we determine slack points, which may be difficult in practice. Du, Parmeter and Racine (2013) show the consistency of the hypothesis test using $D(\hat{p})$ as the test statistic, which implies that the bootstrap version is consistent. In the constrained density setting, the test consists of two steps:

i) If the true density f satisfies the shape constraints, then as $n \rightarrow \infty$,

$$P\{D(\hat{p}) \leq n\epsilon\} \rightarrow 1, \quad \text{for all } \epsilon > 0.$$

ii) If the true function f does not satisfy the shape constraints on \mathcal{J} , then

$$\lim_{\epsilon \rightarrow 0} \liminf_{n \rightarrow \infty} P\{D(\hat{p}) \geq n\epsilon\} = 1.$$

This result has a simple intuitive explanation. If the unconstrained estimator satisfies the constraints, then $D(\hat{p}) = 0$, and clearly there is no need to construct the constrained estimator, because the constraints are most likely true. However, if the constraints are not initially satisfied, then $D(\hat{p})$ can be used to test their validity.

One might consider generalizing the above result to admit different metrics such as, for example, those strictly tailored to probability weights (i.e., $p_i \geq 0$ and $\sum_i p_i = 1$). Although our theory for consistency (Theorem 3) is developed for the Cramér–von Mises statistic, it can instead be developed using a power-divergence metric by following Hall et al. (2001) or by using the L_2 -norm, following Du, Parmeter and Racine (2013). The main difference lies in the algebraic manipulations required for the different metrics. For our theory on the limiting distribution of our test statistic (Theorem 4), we rely on the L_2 -norm, in part because it delivers a tractable solution from the Karush–Kuhn–Tucker conditions. A similar result for, say, the power-divergence metric is possible, though some degree of approximation is still necessary in order to obtain suitable expressions for the weights underlying the corresponding test statistic.

4. Monte Carlo Finite-Sample Performance

In this section, we assess the finite-sample performance of the proposed estimator, and compare it with that of its competitors implemented in currently supported R packages available through CRAN. Note that the proposed estimator is extremely flexible in terms of the type of constraint and the number of simultaneous constraints that can be imposed. For the sake of brevity, we focus on a few test cases, and restrict the group of competitors to the most popular and promising methods of which we are currently aware. The test cases we consider involve two popular constraints, namely the *log-concavity* constraint and the *monotonicity* constraint (i.e., monotonically increasing). Although our approach supports both smooth constrained PDFs and PMFs, we focus on constrained PDF estimation, because of the lack of competing options for smooth constrained PMFs. However, we do provide an illustrative example involving the PMF; the R code for the constrained PDF and PMF estimations is available upon request. The proposed approach can be found in the R package `np` (Hayfield and Racine (2008)), which is available on CRAN. See, in particular, the functions `npuniden.sc()` and `npuniden.boundary()`, which support the constraints monotonically increasing (`constraint="mono.incr"`), decreasing (`constraint="mono.incr"`), convex (`constraint="convex"`), con-

cave (`constraint="concave"`), log-convex (`constraint="log-convex"`), or log-concave (`constraint="log-concave"`), in addition to general inequality constraints placed directly on the density function itself (`constraint="density"` and the upper and lower bound arguments `lb=` and `ub=`). The Cramér-von Mises (`function.distance="TRUE"`) or the L_2 -norm (`function.distance="FALSE"`) can be used to enforce the weights.

For comparison purposes, in the log-concave constraint setting, we compare the proposed approach with those of Cule, Samworth and Stewart (2010), who study a nonsmooth log-concave PDF MLE, and Chen and Samworth (2013), who study the associated smoothed log-concave estimator; these methods can be found in the R package `LogConcDEAD` (Cule, Gramacy and Samworth (2009)) in the functions `mlelcd()` and `dslcd()`. Note that we obtained similar results with the comparable functions in the R package `logcondens` (Dümbgen and Rufibach (2011)), and so do not include these in the analysis below. For an informative overview, see Samworth (2017) for a recent survey of log-concave estimation and its importance in statistical analysis. For comparison purposes, in the monotonically increasing constraint setting, we compare the proposed approach with the monotone rearrangement approach of Birke (2009), which can be found in the R package `Rearrangement` (Graybill et al. (2016)) (see the function `rearrangement()`).

As noted in the introduction, the constrained MLE estimates have a rather nonstandard $n^{-1/3}$ rate of convergence, compared with the $n^{-2/5}$ rate for the kernel estimator. One strength of the MLE approaches is the ease with which they can handle log-concavity in higher dimensions. From a practical perspective, the kernel approach is limited to perhaps $d = 3$ or $d = 4$ dimensions. These approaches are also free of tuning parameters, whereas the kernel approach requires the selection of bandwidths. In the log-concave constraint simulations that follow, we use cross-validation to select the bandwidths for the proposed kernel-based methods, and we optimize the distance from the unconstrained to the constrained *function*, as discussed previously. However, in order to assess the degree to which being free of tuning parameters matters, we begin by comparing the proposed approach based on *infeasible optimal bandwidths* (which are essentially free of tuning parameters) with *data-driven* smoothing parameter selection. Naturally, the optimal bandwidths present the method in the best possible light, albeit an unrealistic one, which is why we use the *data-driven* bandwidth-based results as a reference in the tables that follow, and not the *infeasible* optimal bandwidth-based results. The difference between using the infeasible optimal versus the feasible data-driven tuning parameter (i.e., the bandwidth) is most apparent in small sample settings (e.g., $n = 100$), though this becomes asymptotically negligible as the sample size increases.

In what follows, we consider a modest number of DGPs and, as noted above, restrict our attention to log-concave and monotonically increasing constraints

(the DGPs are presented in Figure 6). The DGPs and a brief description are as follows:

1. The data are drawn from the standard *smooth* unbounded support $N(0, 1)$ univariate Gaussian distribution ($X \in [-\infty, \infty]$), which is log-concave. We report the results based on the (unknown) optimal bandwidth and the data-driven bandwidth, and compare them those of the nonsmooth MLE estimator and the smooth MLE estimator under the log-concavity constraint (Section S1.1).
2. The data are drawn from a *smooth* unbounded support $N(0, \Sigma)$ bivariate Gaussian distribution ($X \in [-\infty, \infty]^2$), which is log-concave, and we compare the results with those of the nonsmooth and the smooth MLE estimators under the log-concavity constraint (Section S1.2).
3. The data are drawn from a *smooth* left-bounded support univariate exponential distribution ($X \in [0, \infty]$), which is log-concave, and we compare the results with those of the nonsmooth MLE estimator and the smooth MLE estimator under the log-concavity constraint (Section S1.3).
4. The data are drawn from a *smooth* bounded support univariate Beta(3,3) distribution ($X \in [0, 1]$), which is log-concave, and we compare the results with those of the nonsmooth MLE estimator and the smooth MLE estimator under the log-concavity constraint (Section S1.4).
5. The data are drawn from a *nonsmooth* bounded support univariate triangular distribution ($X \in [0, 1]$), which is log-concave but *nonsmooth*, and we compare the results with those of the nonsmooth MLE estimator and the smooth MLE estimator under the log-concavity constraint. Note that we include this DGP in order to gauge its robustness, because it violates assumptions invoked when using kernel smoothing methods (i.e., the continuous differentiability of the density up to some particular order > 2) (Section S1.5).
6. The data are drawn from a *smooth* bounded support univariate uniform distribution ($X \in [0, 1]$), which is log-concave, and we compare the results with those of the nonsmooth MLE estimator and the smooth MLE estimator under the log-concavity constraint (Section S1.6);
7. The data are drawn from a *smooth* bounded support univariate Beta(3,1) distribution ($X \in [0, 1]$), which is monotonically increasing, and we compare the results with those of the monotone-rearrangement estimator under the monotonic increasing constraint (Section S1.7).

Note that in each of these scenarios, we report the mean squared error (MSE, computed as $n^{-1} \sum_{i=1}^n (f(X_i) - \hat{f}(X_i))^2$), where $f(\cdot)$ is the true simulation density

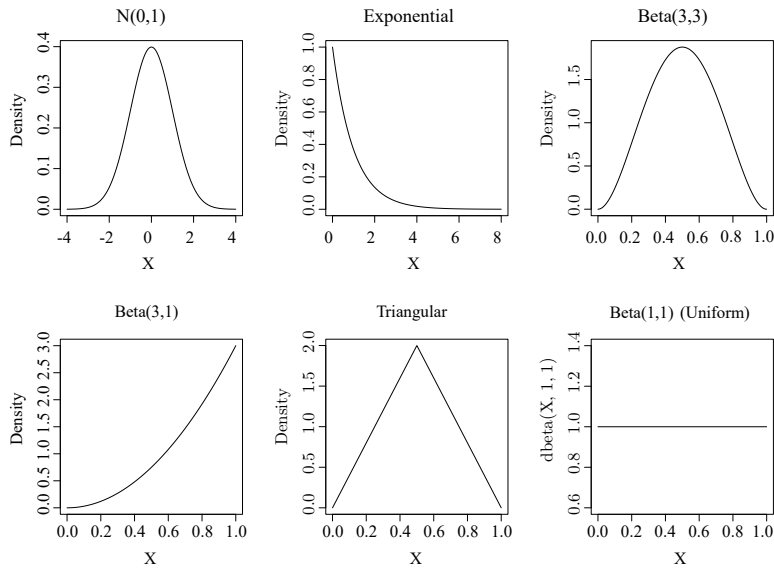


Figure 6. Monte Carlo Densities.

and $\hat{f}(\cdot)$ is an estimate thereof) for the smooth unconstrained version of our estimator (SU) (which is simply the standard kernel density estimation), smooth constrained version of our estimator (SC), nonsmooth MLE estimator (LNS), smooth MLE estimator (LS), and monotone rearrangement estimator (MR). We report the results in both tables (mean/median relative MSE over M Monte Carlo replications) and figures (box plots of the MSE for the M Monte Carlo replications). We present the *relative* MSE values for the mean and median to provide a complete impression of the performance, because the relative mean values may not be robust in the presence of outlying values. Such values occur if, say, data-driven bandwidth selection performs poorly for some fraction of the resamples, and the relative median is naturally less affected by outlying values.

The proposed kernel approach admits known finite boundary points (i.e., the boundary points of $\pm\infty$ have no effect on the estimate), which are used for the exponential (which uses $(a, b) = (0, \infty)$), beta and triangular (each of which use $(a, b) = (0, 1)$) simulations (all other cases use $(a, b) = (-\infty, \infty)$). The peer function `mle1cd()` in the `LogConcDEAD` package does not support known boundary points. Although one might consider modifying the peer function using standard correction methods, it is unclear whether log-concavity is always preserved. Regardless, any such extension of the peer method lies beyond the scope of this study.

In order to meet the page length constraints, the particulars of the Monte Carlo simulations have been moved to the Supplementary Material, which also contains the technical proofs (see Section S1). Briefly, the proposed method is shown to be competitive with its nonsmooth and smooth peers and, most

importantly, provides an extensible and general approach to constraining kernel-based density estimates in a unified framework.

5. Conclusion

We have presented a versatile procedure designed to impose a variety of shape constraints on a *smooth* nonparametric kernel density estimator. We use simulations and real-world data to show that the method can deliver practical and useful estimates of an unknown density, satisfying a range of constraints, and provide the theoretical underpinnings thereof. Additionally, for the constraint of log-concavity, our proposed approach convincingly outperforms popular existing approaches. Furthermore, for the constraint of monotonicity, our approach is competitive with its peers, perhaps even performing somewhat better. However, unlike many of its peers that are tailored for a *single* constraint, our approach is far more flexible and can encompass each of its peers within a unified framework. Moreover, these constraints can be applied to settings involving both continuous and ordered discrete data settings. An R implementation is available on CRAN (see the R package `np` (Hayfield and Racine (2008)), and the functions `npuniden.sc()` and `npuniden.boundary()` contained therein).

There are many exciting and important directions in which the proposed methods can be extended. For example, we can use the insights of Mammen (1991) (in the regression setting) to consider higher-order asymptotic comparisons between the unconstrained and constrained estimators. Given that the constrained density estimator that we propose here is expected to equal the unconstrained estimator *if* the constraints imposed are valid, then for sufficiently large n , these two coincide (to the first order). Hence, one would not expect large sample gains. However, a more nuanced and detailed asymptotic analysis may reveal important higher-order gains that could prove useful for constructing of confidence intervals in small sample settings. Another possible extension is to consider functions supported on a ball, rather than on a hypercube, as considered here. This would require changing existing tools, such as considering kernel functions supported on a ball. Both of these extensions are left to future work.

Supplementary Material

This material contains all of the proofs for the theorems introduced in the paper and the full set of tables and figures for the Monte Carlo simulations.

Acknowledgments

We sincerely thank the editor, the associate editor, and two anonymous reviewers for their helpful comments and suggestions. Du's research was partly supported by the U.S. National Science Foundation grant DMS-1916174.

References

- Aitchison, J. and Aitken, C. G. G. (1976). Multivariate binary discrimination by the kernel method. *Biometrika* **63**, 413–420.
- Bagnoli, M. and Bergstrom, T. (2005). Log-concave probability and its applications. *Economic Theory* **26**, 445–469.
- Birke, M. (2009). Shape constrained kernel density estimation. *Journal of Statistical Planning and Inference* **139**, 2851–2862.
- Chacón, J. E., Duong, T. and Wand, M. P. (2011). Asymptotics for general multivariate kernel density derivative estimators. *Statistica Sinica* **21**, 807–840.
- Chen, Y. and Samworth, R. (2013). Smoothed log-concave maximum likelihood estimation with applications. *Statistica Sinica* **23**, 1373–1398.
- Chernozhukov, V., Fernandez-Val, I. and Galichon, A. (2009). Improving point and interval estimators of monotone functions by rearrangement. *Biometrika* **96**, 559–575.
- Chetverikov, D., Santos, A. and Shaikh, A. M. (2018). The econometrics of shape restrictions. *Annual Review of Economics* **10**, 31–63.
- Cressie, N. A. C. and Read, T. R. C. . (1984). Multinomial goodness-of-fit tests. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **46**, 440–464.
- Cule, M., Gramacy, R. and Samworth, R. (2009). LogConcDEAD: An R package for maximum likelihood estimation of a multivariate log-concave density. *Journal of Statistical Software* **29**, 1–20.
- Cule, M., Samworth, R. and Stewart, M. (2010). Maximum likelihood estimation of a multi-dimensional log-concave density. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72**, 545–607.
- Dette, H., and Pilz, K. F. (2006). A comparative study of monotone nonparametric kernel estimates. *Journal of Statistical Computation and Simulation* **76**, 41–56.
- Du, P., Parmeter C. F. and Racine, J. S. (2013). Nonparametric kernel regression with multiple predictors and multiple shape constraints. *Statistica Sinica* **23**, 1343–1372.
- Dümbgen, L. and Rufibach, K. (2009). Maximum likelihood estimation of a log-concave density and its distribution function: Basic properties and uniform consistency. *Bernoulli* **15**, 40–68.
- Dümbgen, L., and Rufibach, K. (2011). logcondens: Computations related to univariate log-concave density estimation. *Journal of Statistical Software* **39**, 1–28. Web: <http://www.jstatsoft.org/v39/i06/>.
- Fan, J., Zhang, C. and Zhang, J. (2001). Generalized likelihood ratio statistics and wilks phenomenon. *The Annals of Statistics* **29**, 153–193.
- Feng, O. Y., Guntuboyina, A., Kim, A. K. H. and Samworth, R. J. (2021). Adaptation in multivariate log-concave density estimation. *The Annals of Statistics* **49**, 129–153.
- Graybill, W., Chen, M., Chernozhukov, V., Fernandez-Val, I. and Galichon, A. (2016). *Rearrangement: Monotonize oint and interval functional estimates by rearrangement* (version 2.1). Web: <https://CRAN.R-project.org/package=Rearrangement>.
- Grenander, U. (1956). On the theory of mortality measurement. *Scandinavian Actuarial Journal*, 125–153.
- Groeneboom, P. and Jongbloed, G. (2014). *Nonparametric Estimation Under Shape Constraints*. Cambridge University Press.

- Groeneboom, P. and Jongbloed, G. (2018). Some developments in the theory of shape constrained inference. *Statistical Science* **33**, 473–492.
- Hall, P. and Huang, L. S. (2002). Unimodal density estimation using kernel methods. *Statistica Sinica* **12**, 965–990.
- Hall, P., Huang, H., Gifford, J. and Gijbels, I. (2001). Nonparametric estimation of hazard rate under the constraint of monotonicity. *Journal of Computational and Graphical Statistics* **10**, 592–614.
- Hall, P., and Kang, K.-H. (2005). Unimodal kernel density estimation by data sharpening. *Statistica Sinica* **15**, 73–98.
- Hall, P., and Presnell, B. (1999). Density estimation under constraints. *Journal of Computational and Graphical Statistics* **8**, 259–277.
- Hardy, G. H., Littlewood, J. E. and Pólya, G. (1952). *Inequalities*. Cambridge University Press.
- Hausman, J., Hall, B. H. and Griliches, Z. (1984). Econometric models for count data with an application of the Patents-R&D relationship. *Econometrica* **52**, 909–938.
- Hayfield, T. and Racine, J. S. (2008). Nonparametric econometrics: The np package. *Journal of Statistical Software* **27**, 1–32. Web: <http://www.jstatsoft.org/v27/i05/>.
- Horowitz, J. L. and Lee, S. (2017). Nonparametric estimation and inference under shape restrictions. *Journal of Econometrics* **201**, 108–126.
- Koenker, R. and Mizera, I. (2010). Quasi-concave density estimation. *The Annals of Statistics*, 2998–3027.
- Koenker, R. and Mizera, I. (2018). Shape constrained density estimation via penalized Rényi divergence. *Statistical Science* **33**, 510–526.
- Li, Q., and Racine, J. S. (2003). Nonparametric estimation of distributions with categorical and continuous data. *Journal of Multivariate Analysis* **86**, 266–292.
- Li, Z., Liu, G. and Li, Q. (2017). Nonparametric KNN estimation with monotone constraints. *Econometric Reviews* **36**, 988–1006.
- Lok, T. M. and Tabri, R. V. (2021). An improved bootstrap test for restricted stochastic dominance. *Journal of Econometrics* **224**, 371–393.
- Mammen, E. (1991). Estimating a smooth monotone regression function. *The Annals of Statistics* **19**, 724–740.
- Meyer, M. C. (2008). Inference using shape-restricted regression splines. *The Annals of Applied Statistics* **2**, 1013–1033.
- Meyer, M. C., and Habtzghi, D. (2011). Nonparametric estimation of density and hazard rate functions with shape restrictions. *Journal of Nonparametric Statistics* **23**, 455–470.
- Meyer, M. C. and Woodroffe, M. (2004). Consistent maximum likelihood estimation of a unimodal density using shape restrictions. *Canadian Journal of Statistics* **32**, 85–100.
- Meyer-ter-Vehn, M., Smith, L. and Bognar, K. (2017). A conversational war of attrition. *The Review of Economic Studies* **85**, 1897–1935.
- Ouyang, D., Li, Q. and Racine, J. S. (2006). Cross-validation and the estimation of probability distributions with categorical data. *Journal of Nonparametric Statistics* **18**, 69–100.
- Parzen, E. (1962). On estimation of a probability density function and mode. *The Annals of Mathematical Statistics* **33**, 1065–1076.
- Prakasa Rao, B. L. S. (1969). Estimation of a unimodal density. *Sankhya Series A* **31**, 23–36.
- Racine, J., Li, Q. and Yan, K. X. (2020). Kernel smoothed probability mass functions for ordered datatypes. *Journal of Nonparametric Statistics* **32**, 563–586.
- Rathke, F. and Schnörr, C. (2019). Fast multivariate log-concave density estimation. *Computational Statistics & Data Analysis* **140**, 41–58.

- Rosenblatt, M. (1956). Remarks on some nonparametric estimates of a density function. *The Annals of Mathematical Statistics* **27**, 832–837.
- Samworth, R. (2017). Recent progress in log-concave density estimation. *arXiv:1709.03154*.
- Samworth, R. and Sen, B. (2018). Editorial: Special issue on ‘nonparametric inference under shape constraints’. *Statistical Science* **33**, 469–472.
- Tan, G. and Zhou, J. (2020). The effects of competition and entry in multi-sided markets. *The Review of Economic Studies* **88**, 1002–1030.
- Turlach, B. A. and Fortran A. W. (2019). *quadprog: Functions to Solve Quadratic Programming Problems* (version 1.5-8). Web: <https://CRAN.R-project.org/package=quadprog>.
- Walther, G. (2009). Inference and modeling with log-concave distributions. *Statistical Science* **24**, 319–327.
- Wolters, M. A. (2018). *scdensity: Shape-Constrained Kernel Density Estimation* (version 1.0.2.). Web: <https://CRAN.R-project.org/package=scdensity>.
- Wolters, M. A. and Braun, W. J. (2018a). A practical implementation of weighted kernel density estimation for handling shape constraints. *Stat* **7**, e202.
- Wolters, M. A. and Braun, W. J. (2018b). Enforcing shape constraints on a probability density estimate using an additive adjustment curve. *Communications in Statistics - Simulation and Computation* **47**, 672–691.
- Woodroffe, M., and Sun, J. (1993). A penalized maximum likelihood estimate of $f(0+)$ when f is non-increasing. *Statistica Sinica* **3**, 501–515.

Pang Du

Department of Statistics, Virginia Tech, Blacksburg, VA 24061, USA.

E-mail: pangdu@vt.edu

Christopher F. Parmeter

Miami Herbert Business School, University of Miami, Coral Gables, FL 33124, USA.

E-mail: c.parmeter@miami.edu

Jeffrey S. Racine

Department of Economics, McMaster University, Hamilton, Ontario L8S 4L8, Canada.

E-mail: racinej@mcmaster.ca

(Received March 2021; accepted May 2022)