

# Modality-Constrained Density Estimation via Deformable Templates

## Abstract

Estimation of a probability density function (pdf) from its samples, while satisfying certain shape constraints, is an important problem that lacks coverage in the literature. This paper introduces a novel geometric, deformable template constrained density estimator (`dtcode`) for estimating pdfs constrained to have a given number of modes. Our approach explores the space of thus-constrained pdfs using the set of shape-preserving transformations: an arbitrary template from the given shape class is transformed via a shape-preserving transformation to obtain the final optimal estimate. The search for this optimal transformation, under the maximum-likelihood criterion, is performed by mapping transformations to the tangent space of a Hilbert sphere, where they are effectively linearized, and can be expressed using an orthogonal basis. This framework is first applied to (univariate) unconditional densities and then extended to conditional densities. We provide asymptotic convergence rates for `dtcode`, and an application of the framework to the speed distributions for different traffic flows on Californian highways. The supplementary materials for our paper can be found online.

*Keywords:* Density estimation, Modality constraints, Shape constraints, Sieve estimation, Deformation group, Conditional densities,

# 1 Introduction

The estimation of a probability density function (pdf) from its samples is a fundamental problem in statistics, with a multitude of applications in different fields. A subproblem, involving estimation of a pdf given some prior knowledge about the shape of this pdf, is also an important problem. In practice, the prior knowledge stems from a scientific understanding of the underlying process. It is therefore important that the estimate be consistent with the prior shape knowledge in order for it to be interpretable and practically useful as an analytical tool. While a great deal of past research has gone into *shape-constrained* density estimation, these papers have dealt with very specific shape constraints, including log-concavity, monotonicity, and unimodality; there is little to no literature on optimization-based estimation of pdfs with multimodal shape constraints.

The earliest estimate for a unimodal density was given by Grenander (1956), who showed that a particular, natural class of estimators for unimodal densities is not consistent, and presented a modification that is consistent. Over the last several decades, a large number of papers have been written analyzing the properties of the *Grenander estimator*, e.g. Rao (1969); Izenman (1991) and its modifications Birge (1997). An estimator using a maximum likelihood approach was developed by Wegman (1970). The earlier papers assumed knowledge of the position and value of the mode, and applied monotonic estimators over subintervals on either side of it; later papers, for example Meyer (2001); Bickel and Fan (1996), include an additional mode-estimation step. Bayesian methods have also been developed Brunner and Lo (1989). Turnbull and Ghosh (2014), in addition to describing an estimator that uses Bernstein polynomials with the weights chosen to satisfy the unimodality constraint, also provide a useful summary of recent results on unimodal density estimation.

1 The obvious extension to multimodality constraints is of great practical importance because  
2 multimodal densities occur abundantly in nature; in particular, many biological processes are  
3 expected to show a known multimodal structure. For example, the DNA methylation profile in  
4 humans shows a bimodal structure corresponding to hypomethylated and hypermethylated re-  
5 gions: see Harris et al. (2010) and references therein; while the rate of nucleotide substitutions  
6 in DNA sequence (in non-CG-nucleotides) shows a trimodal density corresponding to acceler-  
7 ated, conserved, and neutral substitution rates: see Pollard et al. (2009) and references therein.  
8 In industrial and electrical engineering, household electricity consumption patterns and traffic  
9 patterns have been known to follow multimodal distributions.

## 10 **1.1 Challenges and Current Literature**

11 The important challenges in shape-constrained estimation are to characterize the set of all density  
12 functions satisfying the desired shape constraints, and to solve the maximum likelihood estima-  
13 tion problem on that space. Shape-constrained estimation problems would seem to encourage  
14 a geometric approach, but the use of geometry in density estimation has in fact been sparse in  
15 the literature: to the best of our knowledge, there is *no current method that can impose a mul-*  
16 *timodality constraint on an estimated density and provide optimality in some way.* However,  
17 multimodality constraints have been studied in the case of function estimation: *e.g.* see the very  
18 recent article by Wheeler et al. (2017) and references therein. Here we summarize the literature  
19 that is most relevant to the problem of density estimation under shape constraints.

20 Hall and Huang (2002) introduced a tilting approach to convert an unconstrained density to  
21 an estimate within a *unimodal shape class*. However, the resultant density estimate often directly

1 contradicts the available data by having zero likelihood even at the data points themselves, and  
2 is thus not appropriate as an exploratory tool. Another paper Cheng et al. (1999) proposes to  
3 start with a template unimodal density and provide a sequence of transformations that when  
4 applied to the template both keep the result unimodal, and “improve” the estimate in some sense.  
5 However, the method is *ad hoc*, and asymptotic convergence of the estimates, although seen  
6 empirically, is not guaranteed. Very recently, Wolters and Braun (2018) introduced a technique  
7 that solves the limitations of the approach in Hall and Huang (2002). Specifically, they provide  
8 an algorithm to find a constrained estimate that is the *nearest* to an unconstrained kernel density  
9 estimate (under the integrated squared error loss function), and that can handle up to *bimodal*  
10 constraints. However, this method provides an estimate that satisfies the shape constraint only  
11 on a prespecified grid in the support, so that the estimate need not lie in the correct shape class,  
12 in principle. Since their construction of the constrained estimate involves smoothing out the  
13 spurious peaks of the initial unconstrained estimate, the resultant shape contains spurious flat  
14 spots, which once again limits the interpretability of the estimate. Finally, this estimate is not  
15 designed to be optimal under any specific criterion. This issue is also present in kernel density  
16 estimators, where one can always choose a bandwidth to ensure a given number of modes, but  
17 the resulting density is not optimal in any sense for a finite sample size.

18 Recently, Dasgupta et al. (ress) introduced a geometric approach for exploring the space of *all*  
19 probability densities in order to perform unconstrained density estimation. In this approach, one  
20 starts with an efficient initial estimate, perhaps from a parametric family, and then transforms  
21 it into the desired optimal density using elements of a diffeomorphism group. The problem  
22 therefore shifts to finding the optimal transformation under the chosen criterion. However, no

1 shape constraints are imposed on the estimated density.

## 2 **1.2 Proposed Formulation and Its Novelty**

3 In the current paper, we take a principled and geometrically-intuitive maximum-likelihood ap-  
4 proach to the problem of modality-constrained density estimation. The primary contribution of  
5 this paper is to construct a framework that can handle any general modality constraint, and can  
6 provide smooth interpretable maximum likelihood estimates within a specified shape class.

7 For this purpose, we develop a novel modification of the geometric approach used by Das-  
8 gupta et al. (ress). The method starts with a *template* density from the desired shape class, and  
9 then deforms it into the optimal estimate from that shape class. We shall call this estimator De-  
10 formable Template Constrained Density Estimator or `dtcode`. The advantages of `dtcode` over  
11 existing methods are as follows.

12 First, while estimation is based on deformation or transformation of an initial template as in  
13 Cheng et al. (1999), we apply only a single transformation rather than a possibly non-convergent  
14 sequence. Coupled with a small number of other parameters, this transformation constitutes a  
15 parametrization of the whole of the shape class of interest.

16 Second, we use a broader notion of shape than previous work: in its simplest form, we con-  
17 strain the pdf to possess a fixed, but arbitrary, number of modes; we also consider more gen-  
18 eral cases in the Supplementary Materials. The shape constraint is fully captured in the initial  
19 template itself. As a result, the subsequent estimation of the transformation is independent of  
20 the constraint, providing much greater stability in practical performance with respect to higher  
21 modality constraints than methods such as Wolters and Braun (2018).

1 Third, we use (penalized) maximum likelihood estimation, which guarantees optimality in  
2 principle, and allows the derivation of asymptotic rates of convergence to the true density.

3 The main difference between the current approach and Dasgupta et al. (ress) is in the choice of  
4 transformations used. Dasgupta et al. (ress) wish to parameterize the set of all positive densities.  
5 As a result they choose a set of transformations that act transitively, *i.e.* any positive density  
6 may be transformed into any other. The necessary transformations take the standard form for a  
7 change of variable: a density is transformed by a warping of its domain:  $p \mapsto (p \circ \gamma)\dot{\gamma}$ , where  $p$   
8 is positive probability density and the warping function  $\gamma$  is a diffeomorphism of the domain, *i.e.*  
9 a one-to-one, differentiable map whose inverse is also differentiable. Here,  $\dot{\gamma}$  is the derivative of  
10  $\gamma$ , that is,  $\dot{\gamma}(t) = \frac{d\gamma(t)}{dt}$  for all  $t$  in the domain of the diffeomorphism.

11 Clearly these transformations are not suitable for our case because transitivity is not com-  
12 patible with preserving the shape of a density, merely its normalization. We therefore propose a  
13 different set of transformations, which preserve both normalization and shape: they take the form  
14  $p \mapsto (p \circ \gamma) / \int (p \circ \gamma) dt$ . The denominator renormalizes the density after the transformation in  
15 the numerator; together they preserve the shape of  $p$ , in a sense that we will now explain.

### 16 **1.3 Overview of the Approach**

17 A precise formulation of the problem is as follows: given independent samples  $X = \{x_i\}, i =$   
18  $1, \dots, n$ , from a pdf  $p_0$ , with a known number  $M > 0$  of well-defined modes, estimate this  
19 density ensuring the presence of  $M$  modes in the solution. In order to do this, we construct a  
20 parameterization of the set of continuous densities with  $M$  modes,  $\mathcal{P}_M$ , as follows.

- 21 • Let the set of densities satisfying the shape constraint be denoted  $\mathcal{P}_M = \{p : [0, 1] \rightarrow \mathbb{R}_+ :$

1  $p(0) = p(1) = 0, p$  has  $M$  interior modes}.

2 • Let the critical points of a pdf  $p \in \mathcal{P}_M$  with  $M$  modes be located at  $\{b_a : a \in \{0, \dots, 2M\}\}$ ,  
3 with  $b_0 = 0$  and  $b_{2M} = 1$ .

4 • We define the height ratio vector  $\lambda$  of  $p$  as the set of ratios of the height of the  $(a + 1)^{\text{th}}$   
5 interior critical point to the height of the first (from the left) mode:  $\lambda = (\lambda_1, \dots, \lambda_{2M-2})$ ,  
6 where  $\lambda_a = p(b_{a+1})/p(b_1)$ . Please look at the top left panel of Figure 2 for an illustration.

7 The height ratio vector for the density  $p_0$  illustrated here is simply  $\lambda = (h_2/h_1, h_3/h_1)$ .

8 • Let the subset of  $\mathcal{P}_M$  with height ratio vector  $\lambda$  be denoted  $\mathcal{P}_{M,\lambda}$ . Note that the space  $\mathcal{P}_M$   
9 is the union  $\bigcup_{\lambda} \mathcal{P}_{M,\lambda}$  of the individual spaces  $\mathcal{P}_{M,\lambda}$  with different values of  $\lambda$ .

10 We then parameterize an arbitrary member of  $\mathcal{P}_M$  by:

- 11 1. a height ratio vector  $\lambda \in \Lambda_M$ , where  $\Lambda_M$  is the set of all such vectors;
- 12 2. a diffeomorphism  $\gamma \in \Gamma$ , where  $\Gamma = \{\gamma : [0, 1] \rightarrow [0, 1] : \dot{\gamma} > 0, \gamma(0) = 0, \gamma(1) = 1\}$  is  
13 the group of diffeomorphisms of  $[0, 1]$ . Notably, the set  $\Gamma$  is a group, *i.e.* it is closed under  
14 composition, has an identity element  $\gamma_{id}(t) = t$ , and each element  $\gamma$  has an inverse  $\gamma^{-1}$ .

15 The pdf represented by a pair  $\lambda$  and  $\gamma$  is then  $p_{\lambda,\gamma} = (\tilde{p}_{\lambda}, \gamma) \in \mathcal{P}_{M,\lambda}$ , where  $\tilde{p}_{\lambda} \in \mathcal{P}_{M,\lambda}$  is an  
16 *a priori* fixed template function in  $\mathcal{P}_{M,\lambda}$ , and  $(\cdot, \gamma)$  denotes the transformation of densities by  
17 elements of  $\Gamma$  mentioned earlier, which has the crucial property that it preserves  $\lambda$ .

18 Using this parameterization, we can construct the log likelihood function

$$L(\lambda, \gamma|X) = \sum_i \log p_{\lambda,\gamma}(x_i), \quad (1)$$

19 and we can use maximum likelihood to estimate  $\lambda$  and  $\gamma$ .

1 We can generalize the method to a larger set of shape classes by defining a shape as a se-  
2 quence of piecewise monotonically increasing, decreasing, and flat intervals that together con-  
3 stitute the entire density function. For example, an “N-shaped” density function is given by the  
4 sequence: *increasing-decreasing-increasing*. For any such sequence, we can construct a template  
5 density in the appropriate shape class, and proceed with estimation as before. The assumption  
6  $p_0(0) = p_0(1) = 0$  can also be relaxed, by considering the height ratios of the two boundaries  
7 as two extra parameters  $\lambda_0$  and  $\lambda_{2M+1}$ . We discuss these ideas in more detail in Section 5 of the  
8 Supplementary Materials and present some simulated examples.

## 9 **2 Geometric Representation of Densities**

10 In this section, we show that the above construction does indeed provide a parameterization of  
11  $\mathcal{P}_M$ , by first showing that  $\Gamma$  is large enough to allow us to reach any element of  $\mathcal{P}_{M,\lambda}$  starting  
12 from a template  $\tilde{p}_\lambda \in \mathcal{P}_{M,\lambda}$ , and then showing how to construct such a template for each height  
13 ratio vector  $\lambda \in \Lambda_M$ .

14 **Theorem 1.** *The set of transformations of the set  $\mathcal{P}_{M,\lambda}$  by the mapping  $\mathcal{P}_{M,\lambda} \times \Gamma \rightarrow \mathcal{P}_{M,\lambda}$ , given*  
15 *by  $(p, \gamma) = \frac{p \circ \gamma}{\int (p \circ \gamma) dt}$  is a group action. Furthermore, this action is transitive and free. That is, for*  
16 *any  $p, \tilde{p} \in \mathcal{P}_{M,\lambda}$ , there exists a unique  $\gamma \in \Gamma$  such that  $p = (\tilde{p}, \gamma)$ .*

17 The proof of the theorem is in the Supplementary Materials. The theorem shows that given a  
18 template  $\tilde{p}_\lambda \in \mathcal{P}_{M,\lambda}$ , we can uniquely represent any other pdf  $p$  with the same height-ratio vector  
19 (*i.e.* also in  $\mathcal{P}_{M,\lambda}$ ) as a transformation of the template, *i.e.* as  $p = (\tilde{p}_\lambda, \gamma)$ . What is more, any pdf  
20 in  $\mathcal{P}_{M,\lambda}$  can serve as a template; it can thus be chosen for convenience’ sake.



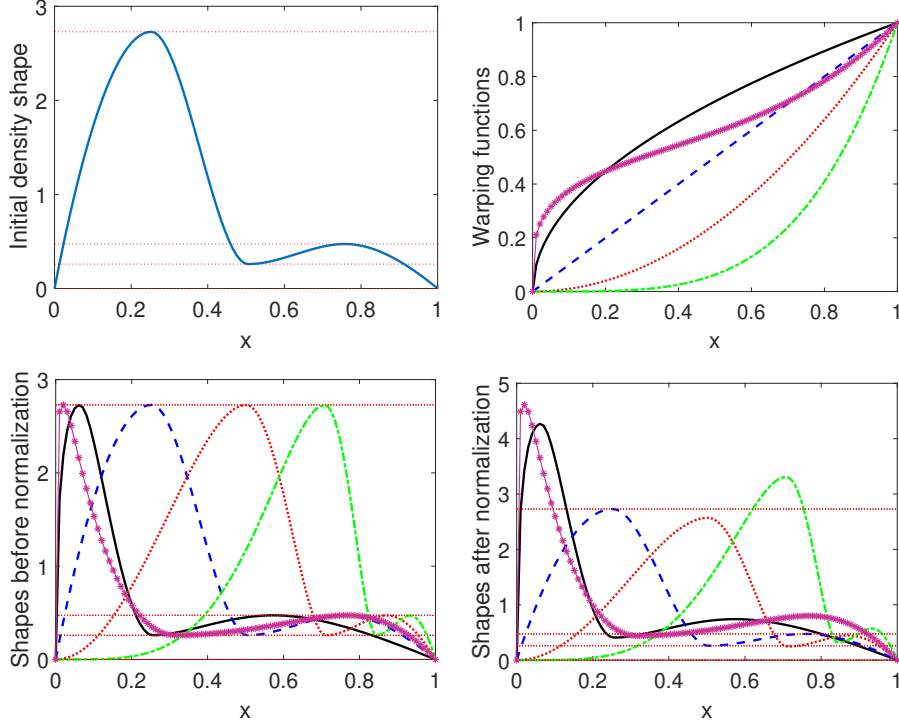


Figure 1: Top left: Initial density. Top right: Different warping functions. Bottom left: Shapes resulting from warping the initial density without renormalization. Bottom right: Resultant warped densities after renormalization.

1      Figure 1 illustrates the height-ratio-vector-preserving effect of the transformations by apply-  
 2      ing several elements of  $\Gamma$  to a pdf in two stages. First, the numerator of the full transformation  
 3      is shown (bottom-left); this stage preserves the heights of all extrema. Second, the pdf is renor-  
 4      malized by dividing by the denominator (bottom-right); this stage changes the heights, but still  
 5      preserves the height-ratio vector.

6      How then do we construct a distinguished template element  $\tilde{p}_\lambda \in \mathcal{P}_{M,\lambda}$ ? First we construct  
 7      an unnormalized function  $g_\lambda$  with  $M$  modes and height ratio vector  $\lambda$ :

- 8      1. Set  $g_\lambda(0) = g_\lambda(1) = 0$ .

- 1     2. Divide the interval  $[0, 1]$  into  $2M$  equal intervals corresponding to the  $M$  modes and  $M - 1$
- 2         interior antimodes, setting  $a_j = j/2M$ ,  $j \in [1, \dots, (2M - 1)]$ , the location of the  $j^{\text{th}}$
- 3         critical point.
- 4     3. Set  $g_\lambda(a_1) = 1$ , and  $g_\lambda(a_j) = \lambda_{j-1}$  for  $j \in [2, \dots, (2M - 1)]$ .
- 5     4. The values of  $g_\lambda$  for all other points are obtained by linear interpolation.

6 We can now define  $\tilde{p}_\lambda = g_\lambda / (\int g_\lambda) \in \mathcal{P}_{M,\lambda}$ .

7 We have thus constructed a representation space  $\Lambda_M \times \Gamma$ , a set of coordinates for  $\mathcal{P}_M$ , where

8  $\Lambda_1 = \{1\}$ , and for  $M > 1$ ,  $\Lambda_M = \{\lambda \in \mathbb{R}_+^{(2M-2)} : \lambda_1 < 1, \lambda_1 < \lambda_2, \lambda_{2j+1} < \lambda_{2j}, \lambda_{2j+1} <$

9  $\lambda_{2j+2}, j = 1, 2, \dots, M - 2\}$ , the conditions arising because the odd indices  $\lambda_1, \lambda_3, \dots, \lambda_{2M-3}$

10 correspond to antimodes, while the rest correspond to modes.

11 Figure 2 shows a simple example to illustrate this representation. The top left panel is a

12 density that has  $M = 2$  modes with critical points located at  $b_i$  and heights  $h_i$ . The top right

13 panel shows the initial template function with  $M = 2$  modes and critical points located at  $a_i$  and

14 heights  $\lambda_i = h_i/h_1$ . The bottom left panel shows the warping function constructed according

15 to the description in the proof of theorem 1, while the last panel shows that using this warping

16 function, we recover the original density.

17 So far, we have assumed that the densities are defined on  $[0, 1]$ . When the bounds of the

18 density function are not known, they are estimated from the data  $X$  using the formula  $A =$

19  $\min(X) - \text{sd}(X)/\sqrt{n}$  and  $B = \max(X) + \text{sd}(X)/\sqrt{n}$ , where  $A$  and  $B$  are the lower and upper

20 bounds respectively,  $\text{sd}(X)$  is the standard deviation of the observations, and  $n$  is the number of

21 observations; these estimates are taken from Turnbull and Ghosh (2014). The data are then scaled

22 to the unit interval,  $z_i = (x_i - A)/(B - A)$ , before proceeding with the rest of the estimation.

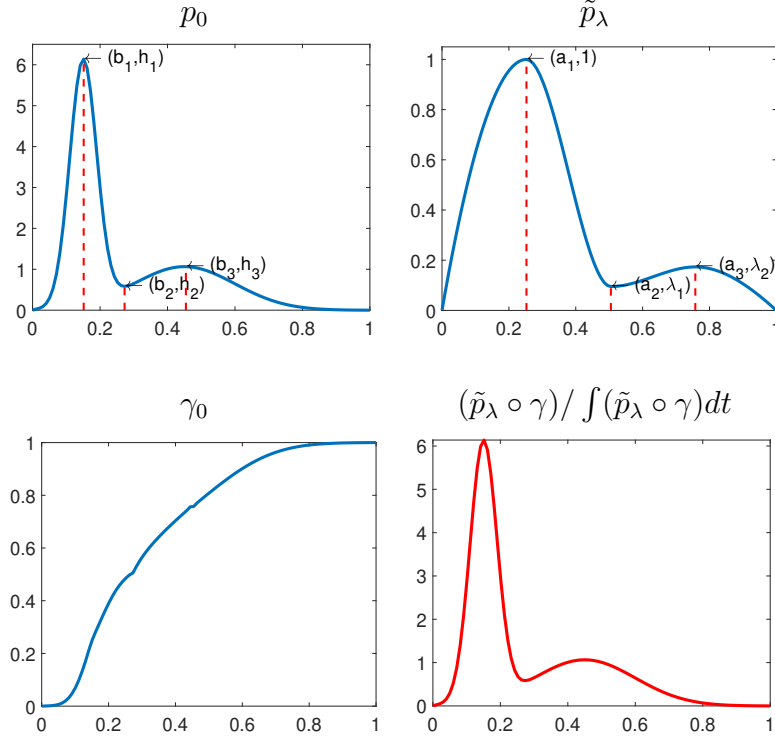


Figure 2: Top left: The original density. Top right: Initial template. Bottom left: The  $\gamma_0$  transforming the template to original shape. Bottom right: Reconstructed density.

1        The framework readily extends to the situation where the true density has a general connected  
 2 support  $\mathcal{D}$  by generalizing from  $\Gamma$  to  $\Gamma^* = \{\gamma : \mathcal{D} \rightarrow \mathcal{D}, \dot{\gamma} > 0, \gamma \text{ is boundary preserving}\}$ . For  
 3 example, if the support of the true density is the entire real line then we can set  $\mathcal{D} = \mathbb{R} \cup \{\pm\infty\}$ .  
 4 However, from a practical standpoint, it is often beneficial to assume that the true density has  
 5 compact rather than infinite support. Our experiments corroborate the findings in Wahba (1981),  
 6 that it is preferable for the true density to have compact support and then to scale the data to  
 7 the unit interval for density estimation. Thus, for the rest of the paper, we always assume that  
 8  $\mathcal{D} = [0, 1]$ , and that the true density has its support on the unit interval.

## 3 Parameter Estimation

Having established a parameterization of the set of shape-constrained densities of interest, the next step is derive a procedure for estimating these parameters from data and specify the pdf estimator `dtcode`. We will use a maximum-likelihood framework, for which we must first specify the log-likelihood function and then solve the optimization problem for  $\lambda$  and  $\gamma$ . The optimization over  $\Gamma$  presents particular difficulties regardless of the likelihood function, and so we first describe how we deal with these.

### 3.1 Finite-Dimensional Representation of Warping Functions

In solving an optimization problem on  $\Gamma$ , we face two challenges. First,  $\Gamma$  is a nonlinear manifold, *i.e.* it is not a vector space; and second, it is infinite-dimensional. We handle the nonlinearity by forming a map from  $\Gamma$  to a vector space. (This vector space happens to be the space tangent to the unit Hilbert sphere  $\mathbb{S}_\infty$  as explained below.) We tackle infinite dimensionality by restricting to a finite-dimensional subspace of this vector space. Together, these two steps are equivalent to finding an increasing family of finite-dimensional subsets  $\Gamma^J \subset \Gamma$  that can be flattened into vector spaces. This then allows us to represent any element  $\gamma \in \Gamma^J$  using a finite orthogonal basis. Once we have a finite-dimensional representation of  $\gamma$ , we can optimize over this representation using standard techniques.

To flatten  $\Gamma$  locally, we define a function  $q_\gamma : [0, 1] \rightarrow \mathbb{R}$ ,  $q_\gamma(t) = \sqrt{\dot{\gamma}(t)}$ , termed the *square-root slope function* (SRSF) of  $\gamma \in \Gamma$ . (For a discussion on SRSFs of general functions, please refer to Chapter 4 of Srivastava and Klassen (2016)). Note that we can reconstruct  $\gamma$  from  $q_\gamma$  using  $\gamma(t) = \int_0^t q_\gamma^2(s) ds$ . In particular, since  $\|q_\gamma\|^2 = \int_0^1 q_\gamma(t)^2 dt = \int_0^1 \dot{\gamma}(t) dt = \gamma(1) - \gamma(0) = 1$ ,

1 we see that  $q_\gamma \in \mathbb{S}_\infty$ , where the unit Hilbert sphere  $\mathbb{S}_\infty$  is defined by  $\mathbb{S}_\infty \subset \mathbb{L}^2 = \{q : [0, 1] \rightarrow$   
2  $\mathbb{R} : \int q^2(t) dt = 1\}$ . We can also see that for any  $q \in \mathbb{S}_\infty$ , there is a  $\gamma_q$  that generates  $q$  given by  
3  $\gamma_q(t) = \int_0^t q^2(s) ds$ .

4 The unit sphere  $\mathbb{S}_\infty$  has known geometry Lang (2012), but is still not a vector space. However,  
5 it can easily be easily flattened into a vector space (locally) due to its constant curvature. A natural  
6 choice for this flattening is a bijective mapping, described next, to the vector space tangent to  
7  $\mathbb{S}_\infty$  at the point  $\mathbf{1}$ , a constant function with value 1. Note that  $\mathbf{1}$  is the SRSF corresponding to  
8  $\gamma = \gamma_{\text{id}}(t) = t$ , *i.e.* the identity, making it a natural choice for the tangent space.) The tangent  
9 space of  $\mathbb{S}_\infty$  at  $\mathbf{1}$  is an infinite-dimensional vector space given by:  $T_1(\mathbb{S}_\infty) = \{v \in \mathbb{L}^2([0, 1], \mathbb{R}) :$   
10  $\int_0^1 v(t)dt = \langle v, \mathbf{1} \rangle = 0\}$ .

11 The bijective mapping between  $\mathbb{S}_\infty$  and  $T_1(\mathbb{S}_\infty)$  is the so-called inverse exponential map:

$$\exp_1^{-1}(q) : \mathbb{S}_\infty \longrightarrow T_1(\mathbb{S}_\infty), \quad v = \exp_1^{-1}(q) = \frac{\theta}{\sin(\theta)}(q - \mathbf{1} \cos(\theta)), \quad (2)$$

12 where  $\theta = \cos^{-1}(\langle \mathbf{1}, q \rangle)$  is the arc-length from  $q$  to  $\mathbf{1}$ .

13 We impose a natural Hilbert structure on  $T_1(\mathbb{S}_\infty)$  using the standard inner product:  $\langle v_1, v_2 \rangle =$   
14  $\int_0^1 v_1(t)v_2(t)dt$ . It is easy to check that, since  $\cos^{-1}(\langle \mathbf{1}, q \rangle) < \pi$ , the norm  $\|v\| = \sqrt{\int_0^1 v(t)^2 dt} =$   
15  $\theta < \pi$  for any  $v = \exp_1^{-1}(q)$ . Thus, the range of the inverse exponential map is not the entire  
16  $T_1(\mathbb{S}_\infty)$ , but a subset  $V = \{v \in T_1(\mathbb{S}_\infty) : \|v\| < \pi\}$ .

17 In order to map points back from the tangent space to the Hilbert sphere, we reverse this  
18 process. This time we use the exponential map:

$$\exp_1(v) : V \rightarrow \mathbb{S}_\infty, \quad \exp_1(v) = \cos(\|v\|)\mathbf{1} + \frac{\sin(\|v\|)}{\|v\|}. \quad (3)$$

19 Finally, we can select any orthogonal basis  $\mathcal{B} = \{b_j, j = 1, 2, \dots\}$  of the Hilbert space  
20  $T_1(\mathbb{S}_\infty)$  and express its elements  $v$  by their corresponding coefficients:  $v(t) = \sum_{j=1}^{\infty} c_j b_j(t)$ ,

1 where  $c_j = \langle v, b_j \rangle$ . The elements of such a basis are just functions in  $\mathbb{L}^2([0, 1], \mathbb{R})$  that are or-  
2 thogonal to  $\mathbf{1}$ , that is,  $\langle b_j, \mathbf{1} \rangle = 0$  for all  $j$ . One example is the Fourier basis excluding  $\mathbf{1}$ , but other  
3 bases, such as the cosine basis, splines, and Legendre polynomials, can also be used. Efromovich  
4 (2010) discusses different choices of basis functions and advocates the use of trigonometric bases  
5 for functions with compact support.

Given a basis  $\mathcal{B} = \{b_j, j = 1, 2, \dots\}$ , one can define an infinite-dimensional space of coeffi-  
cients  $\mathcal{C} = \{c = (c_1, c_2, \dots) : \sum_{j=1}^{\infty} c_j b_j(t) \in V\}$ . One can then truncate the basis expansion to  
approximate elements of  $V$  using finitely-many coefficients. Suppose one uses  $J$  basis elements  
to approximate the tangent space elements. Then, the approximating space of coefficients will  
be denoted by  $\mathcal{C}^J = \{c \in \mathbb{R}^J \mid \sum_{j=1}^J c_j b_j(t) \in V\}$ . Note that  $\mathcal{C}^J$  is a proper subset of  $\mathbb{R}^J$  since  
it contains only elements satisfying  $\|\sum_{j=1}^J c_j b_j(t)\| < \pi$ . Using these two steps, we specify a  
finite-dimensional, and therefore approximate, representation of the transformation space  $\Gamma$ . We  
define a composite map  $H : \mathcal{C}^J \rightarrow \Gamma$ , as

$$\{c_j\} \in \mathcal{C}^J \xrightarrow{\{b_j\}} v = \sum_{j=1}^J c_j b_j \in V \xrightarrow{\text{exp}_1} q \in \mathbb{S}_{\infty} \rightarrow \gamma(t) = \int_0^t q(s)^2 ds. \quad (4)$$

6 For any  $c \in \mathcal{C}^J$ , let  $\gamma_c$  denote the diffeomorphism  $H(c)$ . For any fixed  $J$ , the set  $H(\mathcal{C}^J)$  forms  
7 a  $J$ -dimensional subset of  $\Gamma$ , denoted by  $\Gamma^J$  henceforth, and we pose the estimation problem on  
8 this subset. As  $J$  goes to infinity, this subset  $\Gamma^J$  converges to the full group  $\Gamma$ .

### 9 **3.2 Joint Estimation of $\lambda$ and $\gamma$**

10 We use a joint maximum likelihood method to estimate the height ratios  $\lambda$  along with the coef-  
11 ficients corresponding to the estimate of  $\gamma$ . The maximum likelihood estimate of the underlying

1 density, given the initial template function  $\tilde{p}_\lambda$ , is

$$\hat{p}(t) = \frac{\tilde{p}_{\hat{\lambda}}(\hat{\gamma}(t))}{\int_0^1 \tilde{p}_{\hat{\lambda}}(\hat{\gamma}(t)) dt}, \quad t \in [0, 1], \quad (5)$$

2 where  $\hat{\gamma} = H(\hat{c})$ , and

$$(\hat{c}, \hat{\lambda}) = \operatorname{argmax}_{c \in \mathcal{C}^J, \lambda \in \Lambda_M} \left( \sum_{i=1}^n \left[ \log \left( \tilde{p}_\lambda(\gamma_c(x_i)) / \int_0^1 \tilde{p}_\lambda(\gamma_c(t)) dt \right) \right] \right). \quad (6)$$

3 Since this optimization is over a finite-dimensional Euclidean space, any numerical optimiza-  
 4 tion package can be applied here. The objective function (6) is not convex, and so we use the  
 5 Matlab function `fmincon` for optimization. However `fmincon` can produce local solutions,  
 6 and the `GlobalSearch` toolbox often yields better results, albeit at higher computational cost.  
 7 The `GlobalSearch` toolbox is a multistarting algorithm that generates different trial points as  
 8 initial values of the algorithm, and uses the trial point that converges to a local solution with the  
 9 least objective function value. The algorithm and the method of generating these trial starting  
 10 points are described in Ugray et al. (2007). Depending on the computational resource available,  
 11 one can regulate the number of trial points generated, or can simply use the zero vector as a  
 12 natural starting point.

13 The choice of  $J$ , the number of basis elements, is important. Too large  $J$  can result in over-  
 14 fitting and also put computational burden on the optimization algorithm which might get stuck  
 15 in local, suboptimal solutions. We use a penalized version of the likelihood in (6), the standard  
 16 Akaike's Information Criterion (AIC), to choose the optimal number of basis elements.

## 1 **4 Simulation Study**

2 For the numerical implementation of `dtcode`, we use the Fourier basis for the tangent space  
3 representation. We start with 2 basis elements, and increase the number up to a pre-decided  
4 limit; we then choose the result with the best AIC value. We chose AIC as the penalty on the  
5 number of basis elements because experiments suggested that BIC overpenalizes the number of  
6 parameters, causing the estimate to miss the sharper features of the true density. The code for  
7 `dtcode` is available online at [https://github.com/Sutanoy/Shapeconstrained\\_](https://github.com/Sutanoy/Shapeconstrained_DensityEstimation)  
8 `DensityEstimation`.

9 For illustration, we use sample sizes of 100, 500, and 1000. To evaluate the average perfor-  
10 mance, we generate 100 samples (of sample size 100, 500, and 1000 respectively) and evaluate  
11 the mean error and the standard deviation of the errors. For the error function, we considered the  
12 vector  $\mathbb{L}^2$ ,  $\mathbb{L}^1$ , and  $\mathbb{L}^\infty$  norms of the difference between the true density and the density estimate  
13 evaluated on 100 equidistant points across the support.

14 The average computational time for `dtcode` varies from around 20 seconds for a sample of  
15 size 100, to 250 seconds for a sample of size 1000, while optimizing over ten different possible  
16 parameter dimensions using 1000 trial points in the `GlobalSearch` algorithm, on an Intel(R)  
17 Core(TM) i7-3610QM CPU processor laptop.

### 18 **Study 1**

19 We start with two examples with one mode:

- 20 • Example 1: a symmetric unimodal pdf given by  $p_0 = 0.8\mathcal{N}(0, 4) + 0.2\mathcal{N}(0, 0.5)$ .
- 21 • Example 2: a unimodal pdf with contamination, given by  $p_0 = 0.95\mathcal{N}(0, 0.5) + 0.05\mathcal{N}(3, 1)$ .



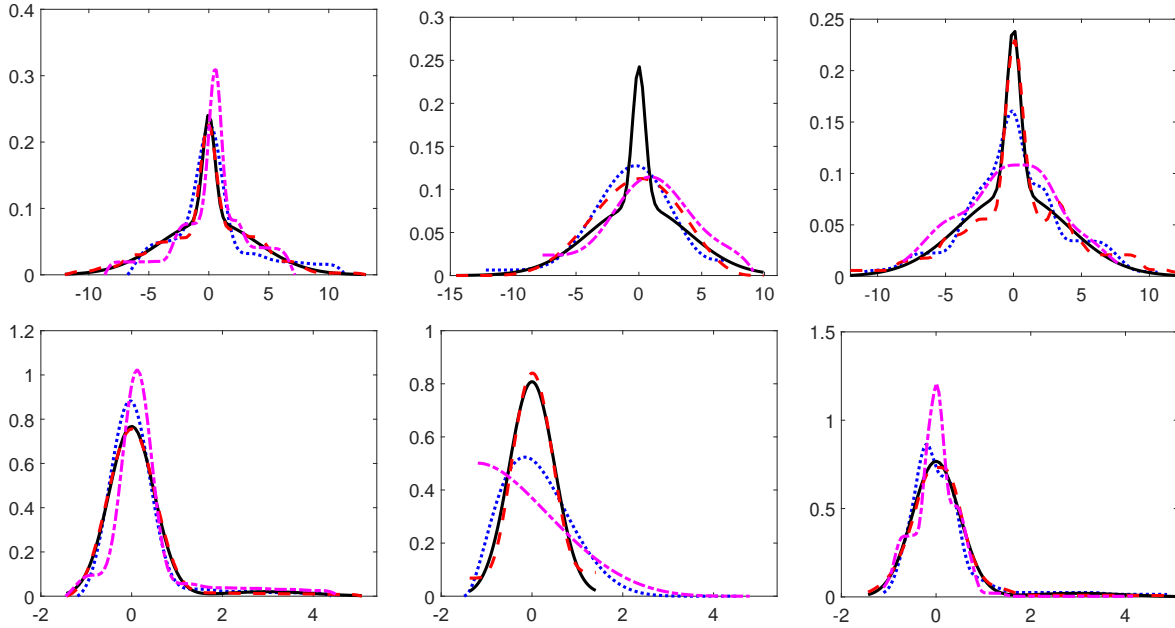


Figure 3: Performance of `dtcode` (left) vs. `umd` (centre) and `scdensity` (right) for examples 1 (top row) and 2 (bottom row) of Study 1, using the  $\mathbb{L}^2$  loss function values calculated for 100 samples of size 100. The true density is shown as a solid line; the estimated density with best performance as a dashed line; with median performance as a dotted line; and with worst performance as a dashed-dotted line.

1 In Figure 3, we use the  $\mathbb{L}^2$  loss function values calculated for 100 samples of size 100 to compare  
 2 the results obtained with `dtcode` (leftmost column) to those obtained using the `umd` package de-  
 3 veloped by Turnbull and Ghosh (2014) (centre column) and the `scdensity` package introduced  
 4 in Wolters and Braun (2018) (rightmost column). The upper and lower rows show examples 1  
 5 and 2. In each plot, the true density is shown as a solid line; the estimated density: with best per-  
 6 formance as a dashed line; with median performance as a dotted line; and with worst performance  
 7 as a dashed-dotted line.

8 In both examples, `dtcode` clearly outperforms `umd` in capturing the sharper features and in

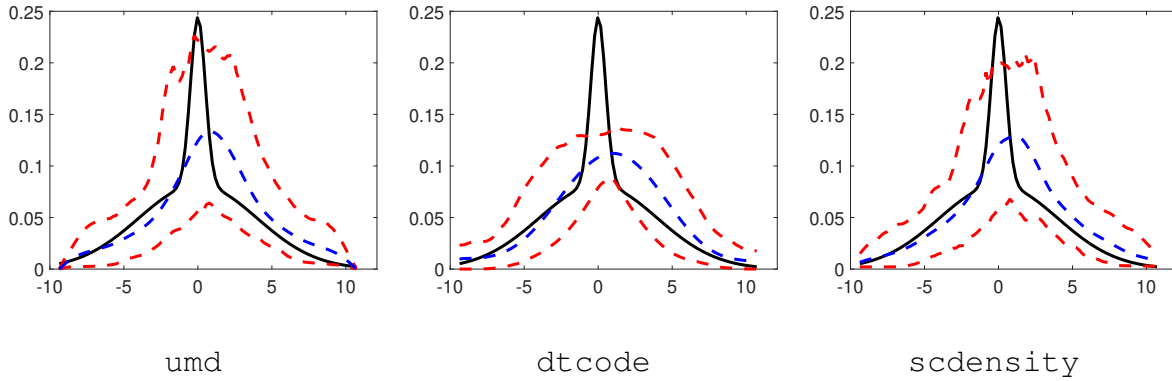


Figure 4: A comparison of the variability of the estimates across different samples for example 1 from Study 1 at sample size 100. The middle dashed line indicates the average estimate across samples, while the upper and lower dashed lines represent the 95<sup>th</sup> and 5<sup>th</sup> quantiles, respectively, of the estimate at the location. The solid line is the true density.

1 stability of performance. On the other hand, the performance of `dtcode` is very similar to  
 2 the performance of `scdensity` for both the examples. For the kurtotic unimodal example 1,  
 3 `scdensity` has a slightly better performance, whereas for the contaminated unimodal density  
 4 estimate 2, `dtcode` is superior. Table 1 in the Supplementary Materials gives a quantitative  
 5 analysis.

6 For the kurtotic unimodal example 1, we also study the pointwise MSE, as shown in Figure 4.  
 7 The figure shows that across all samples, `dtcode` and the package `scdensity` have a similar  
 8 overall performance in capturing the location and the height of the mode.

9 Example 2 is a special case where there are outliers in the data that create the possibility of a  
 10 small peak near the right boundary. These outliers also affect the boundary estimates, and reflect  
 11 a spurious mode in the true density. In this example, we see that `dtcode` is very robust to the  
 12 choice of boundary estimates, replacing the spurious mode with a wide shoulder, as shown in the

1 bottom left panel of Figure 3. For this example, the quantitative performance of `dtcode` is also  
2 superior to the other techniques for all sample sizes, as shown in Table 1 of the Supplementary  
3 Materials.

#### 4 **Study 2**

5 Next, we study a bimodal density: an asymmetric bimodal density given by  $p_0 = 1/3\mathcal{N}(-1, 1) +$   
6  $2/3\mathcal{N}(1, 0.3)$ .

7 We compare the estimation performance of `dtcode` with `scdensity`. Table 2 in the Sup-  
8 plementary Materials presents a quantitative comparison of different loss functions and the like-  
9 lihood for this example at different sample sizes. For all sample sizes, the performance of the  
10 two approaches is very similar with respect to the loss functions. However, there are some clear  
11 advantages to our approach. First, note that the bimodality constraints in Wolters and Braun  
12 (2018) are satisfied only on a pre-specified finite grid. As a result, the final estimate has spurious  
13 modes violating the shape constraint, and thus technically does not belong in the correct shape  
14 class; the ability to violate the constraints probably explains the slightly better  $\mathbb{L}^2$  errors for its  
15 estimates. Secondly, the estimate itself does not enjoy any statistical optimality. The estimate  
16 starts with an unconstrained estimate and obtains the nearest estimate in the correct shape class.  
17 For that purpose, it replaces spurious peaks with flat intervals even though the data might suggest  
18 otherwise.

19 Figure 5 illustrates an example of the performance of `dtcode` in comparison with `scdensity`.

20 The left panel shows the `dtcode` result, while the right panel shows that for `scdensity`.  
21 The 100 observations are also shown along the horizontal axis. The quantitative performance of

1 `scdensity` and `dtcode` (shown in Tables 1 and 2 of the Supplementary Materials) are very  
2 similar at all sample sizes. Further investigation reveals that small differences can mostly be  
3 attributed to the choice of starting point and the actual optimization algorithm used in our ap-  
4 proach, rather than the approach itself. For example, we notice that `scdensity` performs better  
5 if we use an external optimization function to obtain the mode locations rather than the approach  
6 proposed in the original paper and then used in the `scdensity` package. Also, `dtcode` shows  
7 improvement if we choose a more informed starting template shape, such as a kernel density with  
8 hand-tuned bandwidth so that the number of modes is correct.

9 With respect to the shape of the resultant density estimate, however, the `scdensity` es-  
10 timate does not conform to the available data because the constraint is only satisfied on a pre-  
11 specified grid. The right panel of Figure 5 shows the `scdensity` estimate along with the local  
12 maxima indicated by asterisks. The left modal region is replaced by several small bumps, mak-  
13 ing it difficult to distinguish a true mode from the constraint violations. We also note that the  
14 spurious flat shape in the left tail is probably due to the inbuilt optimization code provided. In  
15 comparison, `dtcode` correctly captures the data-sparse region in between the two modal regions  
16 and has exactly two modes.

17 Finally, we emphasize that the shape constraints appear directly in the estimation procedure of  
18 Wolters and Braun (2018). This makes the constrained estimation and the nested search for criti-  
19 cal points increasingly complex as the modality is increased and makes the approach ill-equipped  
20 to handle higher modality constraints. In contrast to `scdensity`, the constraint information in  
21 `dtcode` is captured in the initial template function itself, and the subsequent estimation of the  
22 transformation is free of the modality information, meaning that the approach scales much better

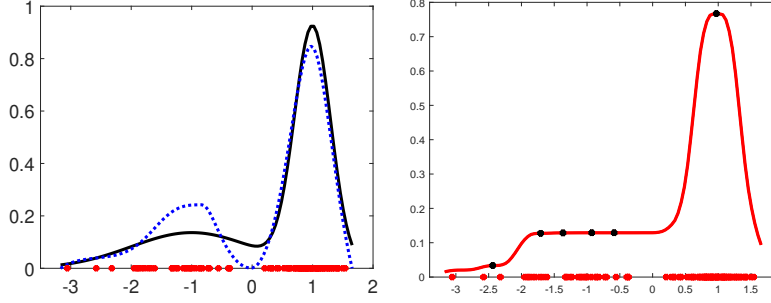


Figure 5: The left panel shows the bimodal example from Study 2, with the true density (solid) and `dtcode` estimate(dotted). The right panel shows the `scdensity` estimate along with local modes. The observations are shown on the  $x$ -axis.

1 to more general modality constraints.

## 2 Study 3

3 As an extension of the previous experiments, we now study performance across a range of uni-  
 4 modal and bimodal examples. We do this by averaging performance over random samples from  
 5 a set of random densities in the appropriate shape family. The true densities themselves are  
 6 generated randomly as follows:

- 7 • Example 1: a unimodal example with random mixing proportions and standard deviations,  
 8 given by  $p_0 = \alpha\mathcal{N}(0, \sigma_1) + (1-\alpha)\mathcal{N}(0, \sigma_2)$ , with  $\alpha \sim U(0, 1)$ ,  $\sigma_1 \sim \max(0.1, \mathcal{N}(0.4, 0.1))$ ,  
 9 and  $\sigma_2 \sim \max(0.1, \mathcal{N}(3, 0.2))$ .
- 10 • Example 2: a bimodal example with random mixing proportions, means, and standard  
 11 deviations, given by  $p_0 = \alpha\mathcal{N}(\mu_1, \sigma_1) + (1 - \alpha)\mathcal{N}(\mu_2, \sigma_2)$ , with  $\alpha \sim U(0, 1)$ ,  $\sigma_1 \sim$   
 12  $\max(0.1, \mathcal{N}(0.75, 0.2))$ ,  $\mu_1 \sim \mathcal{N}(-1, 0.2)$ ,  $\mu_2 \sim \mathcal{N}(0.1, 0.2)$ , and  $\sigma_2 \sim \max(0.1, \mathcal{N}(0.5, 0.2))$ .

13 Figure 6 shows boxplots of the  $\mathbb{L}^2$  norms of the estimation errors for example 1 (top row)

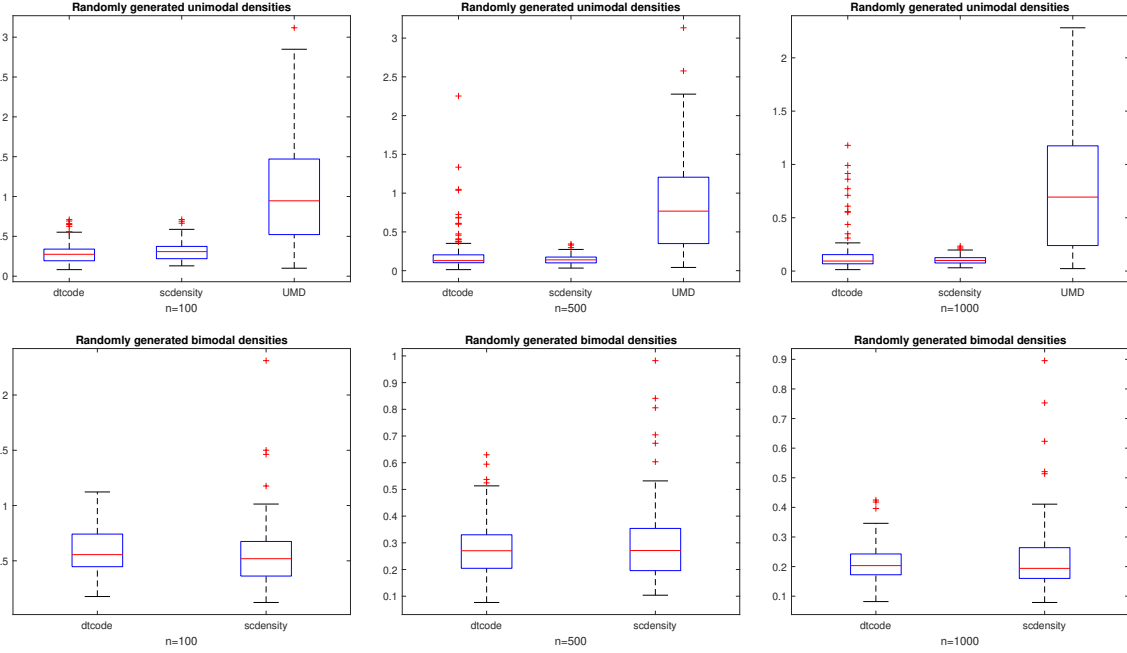


Figure 6: Boxplots of the  $\mathbb{L}^2$  norms of the errors of the density estimates found using `dtcode`, `scdensity`, and `umd`, for randomly sampled true densities and different sample sizes. The top row shows the results for example 1; the bottom row those for example 2.

1 and example 2 (bottom row), for three sample sizes. We notice that the `dtcode` estimate is  
 2 comparable to the `scdensity` estimate at all sample sizes under this measure. Again, as above,  
 3 our approach is better in terms of the desired shape constraint.

#### 4 Study 4

5 Next, we study pdfs with three and four modes respectively:

- 6 • Example 1: an asymmetric trimodal density with one mode well separated from the other  
 7 two, given by  $p_0 = 1/3\mathcal{N}(-1, 0.25) + 1/3\mathcal{N}(0, 0.25) + 1/3\mathcal{N}(2, 0.3)$ .
- 8 • Example 2: a four-modal density, given by  $p_0 = 0.25\mathcal{N}(-4, 0.5) + 0.25\mathcal{N}(-2, 0.5) +$

1  $0.4\mathcal{N}(2, 1) + 0.1\mathcal{N}(5, 0.25)$ .

2 In Figure 7, we use the  $\mathbb{L}^2$  norm of the errors calculated for 100 samples of size 100 to study  
3 the results obtained with `dtcode` on examples 1 (left column) and 2 (right column). The top  
4 row shows plots of the true density as a solid line; the estimated density: with best performance  
5 as a dashed line; with median performance as a dotted line; and with worst performance as a  
6 dashed-dotted line. The bottom row shows boxplots of the  $\mathbb{L}^2$  norms of the errors for different  
7 samples sizes. The results show that the performance improves with increasing sample size, in  
8 both size and spread of error. Note that we do not compare the `dtcode` results to those of other  
9 methods because there is no other method that can handle  $M = 3$  or higher.

## 10 **5 Application To Electricity Consumption Data**

11 Quantification and detection of patterns in electricity consumption curves across households, lo-  
12 cations, and seasons, is crucial for planning and forecasting, as discussed in Cordova et al. (2018)  
13 and Kwac et al. (2014), among others. Deployment of advanced monitoring systems, including  
14 smart meters and synchrophasors, in power distribution networks has created a new paradigm  
15 for observing and managing the electric grid, leading to an abundance of consumption data with  
16 different levels of granularity. The City of Tallahassee, the capital of Florida, has a Meter Data  
17 Management System (MDMS) that stores electricity consumption (kWh) readings from every  
18 customer in the city for billing purposes and further analysis. We look at the daily electricity  
19 consumption profiles of a randomly chosen de-identified single household in Tallahassee. The  
20 dataset was obtained with a Non-Disclosure Agreement with the City of Tallahassee.

21 The daily consumption patterns show high variability, depending on day of the week, season,

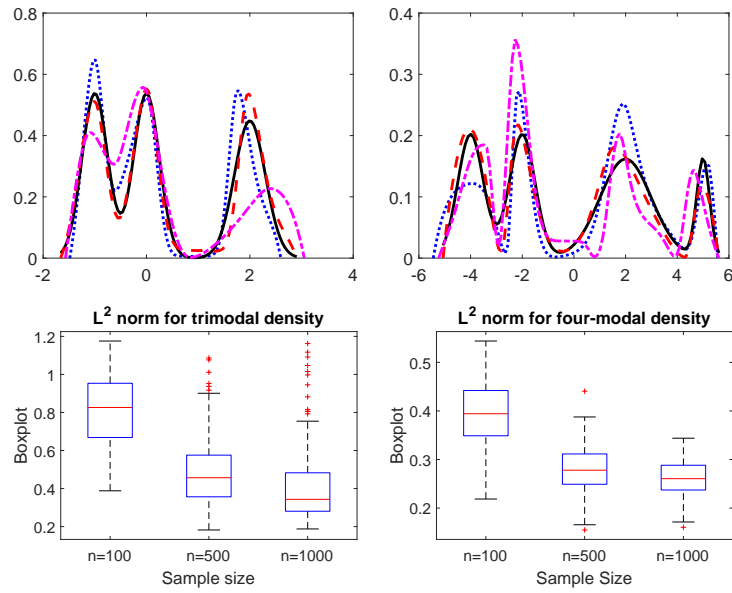


Figure 7: Results of the `dtcode` method on examples 1 (left column) and 2 (right column) of Study 4. The top row shows plots of the true density as a solid line; the estimated density: with best performance as a dashed line; with median performance as a dotted line; and with worst performance as a dashed-dotted line; with performance measure using the  $L^2$  norms of the errors. The bottom row shows boxplots of the  $L^2$  norms of the errors for different samples sizes.



1 and other extraneous factors, even for this single household. We look at the electricity consump-  
 2 tion values at different times in a particular day, for four different days, in order to estimate  
 3 the daily distribution of electricity consumption. However, one can split the daily consumption  
 4 profiles into two interpretable clusters: consumption values when the household members are at  
 5 home versus consumption values when the households are not at home. This suggests that a two-  
 6 mode constrained density estimation would lend interpretability to the density estimates, which  
 can otherwise be very noisy.

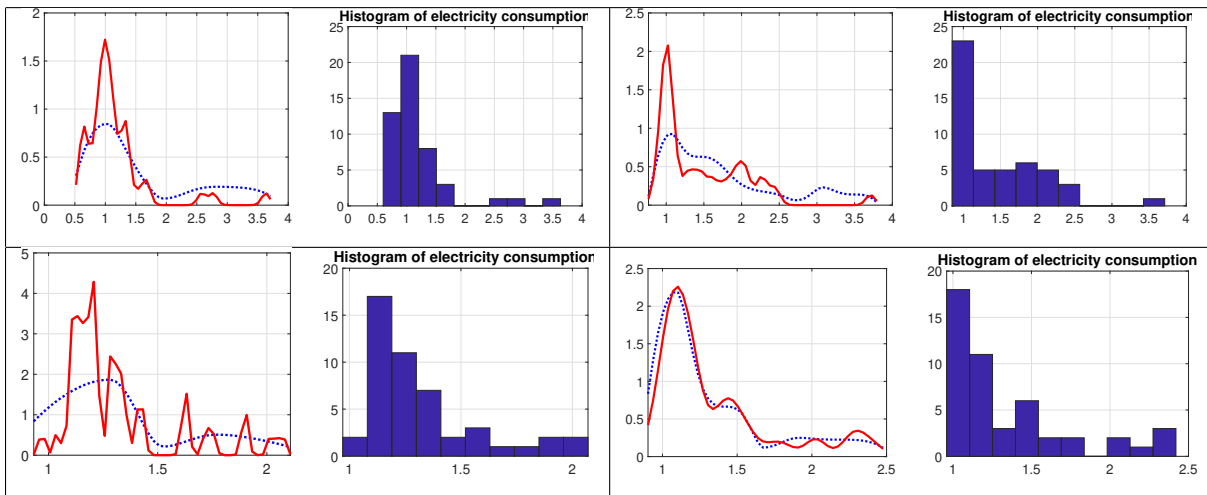


Figure 8: Estimated densities of electricity consumption using the warped approach (dotted) and a kernel (ucv) approach (solid), and the associated histograms of consumption data for four different days.

7

8 Figure 8 shows the density estimates and the corresponding histograms of electricity con-  
 9 sumption for four different days at the randomly chosen household. As expected, in most cases,  
 10 an unconstrained estimate is too bumpy and uninterpretable. The shape constraint, however, re-  
 11 sults in much smoother and more interpretable estimates, the two peaks captured by our proposed  
 12 method aligning well with the major peaks in the histograms.

## 6 Extension To Conditional Density Estimation

The proposed framework for modality-constrained density estimation extends naturally to modality-constrained *conditional* density estimation. Consider the following setup. Let  $X$  be a fixed one-dimensional random variable with a positive density on a fixed interval. Let  $Y \sim f_X$ , where  $f_X$  is an unknown conditional density that changes smoothly with  $X$ .

Conditioned on  $X$ ,  $Y$  is assumed to have a univariate, continuous distribution with support on the interval  $[A, B]$ , with  $M$  modes in the interior of  $[A, B]$ , and  $f_X(A) = f_X(B) = 0$ . We observe the pairs  $\{(Y_i, X_i)\}, i = 1, \dots, n$ , and are interested in recovering the conditional density  $f_X$  at a particular location of  $X$ , henceforth referred to as  $x_0$ . The estimation is again initialized with an  $M$ -modal template function  $\tilde{p}_\lambda$ . However, since  $f_X$  varies smoothly with  $X$ , we assign more importance to observations closer to the location  $x_0$  than to observations further away, and hence we perform weighted maximum likelihood estimation to find the necessary parameters:

$$(\hat{c}, \hat{\lambda}) = \underset{c \in \mathcal{C}^J, \lambda \in \Lambda_M}{\operatorname{argmax}} \left( \sum_{i=1}^n \left[ \log \left( \tilde{p}_\lambda(\gamma_c(x_i)) / \int_0^1 \tilde{p}_\lambda(\gamma_c(t)) dt \right) \right] W_{x_0, i} \right) \quad (7)$$

where  $W_{x_0, i}$  is the localized weight associated with the  $i^{\text{th}}$  observation, calculated according to:

$$W_{x_0, i} = \frac{\mathcal{N}(\|X_i - x_0\|_2 / h(x_0); 0, 1)}{\sum_{j=1}^n \mathcal{N}(\|X_j - x_0\|_2 / h(x_0); 0, 1)}. \quad (8)$$

Here  $\mathcal{N}(\cdot, 0, 1)$  is the standard normal pdf and  $h(x_0)$  is the parameter that controls the relative weights associated with the observations. However, weights defined in this way result in higher bias because information is being borrowed from all observations. To mitigate this, as discussed in an example in Bashtannyk and Hyndman (2001), we allow only a specified fraction of the observations  $X_i$  to have a positive weight. Note that using too small a fraction will result in unstable estimates and poor practical performance because the effective sample size will be too

1 small. Hence we advocate using the 50% of the observations nearest to the target location for  
2 borrowing information, and then calculating the weights for this smaller sample as before.

3 The parameter  $h(x_0)$  is akin to the bandwidth parameter associated with traditional kernel  
4 methods for density estimation, for the predictors  $X$ . A very large value of  $h(x_0)$  distributes  
5 approximately equal weight over all observations, whereas a very small value considers only the  
6 observations in a small neighbourhood around  $x_0$ . The value of  $h(x_0)$  can be chosen via any stan-  
7 dard cross-validation-based bandwidth selection method. In our experiments, we use an adaptive  
8 bandwidth selection method to save computation time when the predictors are independent of  
9 each other. It consists of a two-step procedure:

10 1. Compute a standard kernel density estimate  $\hat{K}$  of the predictor space using a fixed band-  
11 width chosen according to any standard criterion. (We simply use the `ksdensity` esti-  
12 mate in MATLAB, which chooses the bandwidth optimal for normal densities.) Let  $h$  be  
13 the fixed bandwidth used.

14 2. Then, set the bandwidth parameter  $h(x_0)$  at location  $x_0$  to be  $h(x_0) = h/\sqrt{\hat{K}(x_0)}$ .

15 The intuition behind this choice is that  $h$  controls the overall smoothing of the predictor space  
16 based on the sample points, while  $\sqrt{\hat{K}(x_0)}$  stretches or shrinks the bandwidth at the particular  
17 location. In a sparse region, increased borrowing of information from other data points is desir-  
18 able in order to reduce the variance of the estimate, whereas in dense regions, reduced borrowing  
19 of information from faraway points reduces the bias of the density estimates. A location from  
20 a sparse region is expected to have a low density estimate, and a location from a dense region  
21 is expected to have a high density estimate. Hence, varying the bandwidth parameter inversely  
22 with the density estimate helps adapt to the sparsity around the point of interest. The choice of

1 the adaptive bandwidth parameter is motivated by the variable bandwidth kernel density estima-  
2 tors discussed in Terrell and Scott (1992), Van Kerm et al. (2003), and Abramson (1982), among  
3 others. We provide a simulation study in the Supplementary Materials.

## 4 **7 Application To Traffic Flow Data**

5 As an application of modality-constrained conditional density estimation, we use the traffic flow  
6 data for Californian highways from the package `hdrcde` in R. The scatterplot shown in Figure 9  
7 shows the distinctly bimodal nature of the speed distribution for traffic flows between 1000 and  
8 1620 vehicles per lane per hour, corresponding to uncongested and congested traffic. This range  
9 of traffic flows has already been studied by Einbeck and Tutz (2006). They note that beyond a  
10 traffic flow of 1620, the regression curves corresponding to uncongested and congested traffic  
11 are no longer distinguishable. So, we consider the speed flow in the above range (772 observa-  
12 tions), and estimate the density of the speed conditional on a flow of 1400. We use a bimodality  
13 constraint on the shape, and our prescribed 50% of the 772 observations. For the tangent space  
14 representation, we use up to 6 basis elements.

15 The middle panel of Figure 9 (solid line) shows the conditional density estimate for flow  
16 = 1400 using `dtcode`. The left mode is at 35.56mph and the right mode is at 59.01mph. Ein-  
17 beck and Tutz (2006) obtain a very similar conditional density estimate. The left mode in their  
18 case is at 32.65mph and the right mode is at 59.18mph. On the other hand, if we find a tradi-  
19 tional conditional density estimate using the NP package, we find several spurious bumps; this  
20 estimate is shown in the middle panel of Figure 9 (dotted line), with a magnified view shown  
21 in the right panel. The superfluous bumps are present in the NP estimate constructed using 772

1 observations (not presented), as well as the estimate constructed using only 50% of the obser-  
 2 vations as in our approach. This results in over-interpreting the tail and consequently a lack of  
 3 interpretability for the modes themselves. Thus constraining the number of modes clearly helps  
 with the interpretability of the resulting density.

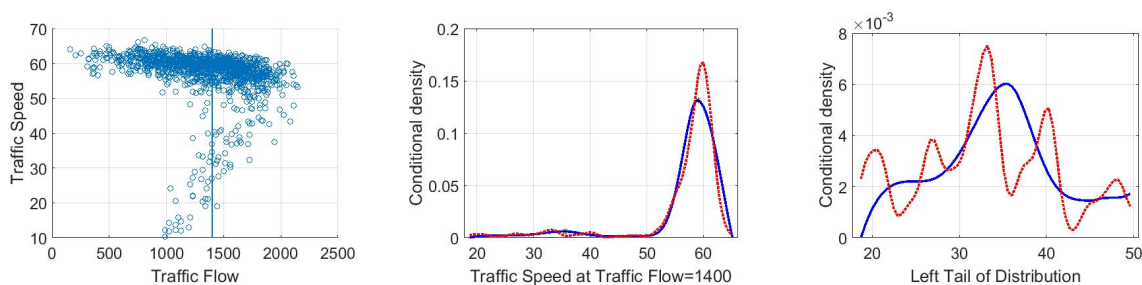


Figure 9: Left: scatterplot of traffic flow data for Californian highways from the `hdrcde` package in R. Centre: traffic speed density at traffic flow 1400 as estimated by `dtcode` (solid line) and the `NP` package (dotted line). Right: a magnified view of the left part of the centre plot.

4

## 5 **8 Discussion**

6 Shape-constrained density estimation is a rich problem area that has a broad range of real-world  
 7 applications, yet has been explored rigorously only in limited cases. Here we have introduced a  
 8 novel framework, using geometric tools, that enables shape-constrained density estimation using  
 9 a different notion of shape than studied previously. In our approach, named `dtcode`, a template  
 10 from the appropriate shape class is deformed using shape-preserving diffeomorphisms of the data  
 11 domain, the optimal deformation being defined by maximum likelihood. The problem is thereby  
 12 reframed as one of optimizing over the diffeomorphism group.

13 The framework is the first in the literature that can perform modality-constrained density

1 estimation for any number of modes. However, from a practical perspective, the performance  
2 suffers somewhat when the constrained shape becomes too complex or if the number of modes  
3 becomes high (greater than 4). This limitation is due to the current choice of numerical techniques  
4 used in optimization over the diffeomorphism group, and because of the choice of basis set used  
5 in estimation.

6 Since this paper primarily focuses on the fundamental framework for `dt_code`, it only lightly  
7 touches upon or leaves out some associated problems. Examples include the choice of the number  
8 of basis elements for the tangent space representation, the choice of the basis itself, estimation  
9 of domain boundaries, and the choice of penalty for penalized estimation. These are all inter-  
10 esting problems in their own right, but space limitations force our focus to only the main ideas.  
11 Nevertheless, we can make some observations.

- 12 • This paper uses AIC as the penalty to select the number of basis elements because, in  
13 comparison, BIC tends to choose an insufficient number of parameters. However, other  
14 model selection techniques can also be investigated.
- 15 • Experiments using a Meyer wavelet basis for the tangent space representation yielded re-  
16 sults similar to those reported in the paper, although the Meyer wavelets seemed to require  
17 more observations than the Fourier basis to obtain satisfactory results. Clearly, one can  
18 choose different bases and conduct a comparative study of performance. Since the support  
19 of the warping functions is compact, we recommend using trigonometric (Fourier and co-  
20 sine) basis for representation. Please refer to Efromovich (2010) and the references therein  
21 for a more detailed discussion on this topic. When the sample size is small, Fourier basis  
22 can result in spurious bumps near the boundaries, which is why wavelets may be a good

1 alternative.

- 2 • Our paper follows Turnbull and Ghosh (2014) in estimating the boundaries, but other  
3 choices can be explored as well.
- 4 • For conditional density estimation, the weights can be defined using any kernel: the Gaus-  
5 sian kernel (and the  $\mathbb{L}^2$  -loss function) was only used as an illustration.

6 An advantage of the proposed framework is that it is easy to extend to conditional density  
7 estimation via a weighted maximum likelihood objective function. One potential future direction  
8 is to apply this framework to situations where a large number of covariates are present. Currently  
9 the bandwidth parameter is chosen adaptively based on a kernel density estimate at the location  
10 of the (scalar) covariate. The framework can be directly extended to a scenario with  $d$  covariates  
11 using a  $d$ -dimensional kernel density estimate at the location of the predictors. Such an esti-  
12 mate would generically suffer from the curse of dimensionality, but seems valid for applications  
13 where only a few of the covariates are relevant. In particular, Wasserman and Lafferty (2006)  
14 have developed a technique to shortlist relevant variables and to find corresponding bandwidth  
15 parameters. Using these bandwidth parameters, one can redefine the weights and then perform  
16 weighted likelihood maximization as before to produce a conditional density estimate.

17 In conclusion, we have developed a framework for incorporating general modality constraints  
18 into a density estimation procedure, while showing very competitive performance on shape con-  
19 straints already studied in literature. In applications where the data shows modality constraints,  
20 the proposed framework will provide accurate and interpretable density estimates that fully re-  
21 spect the constraints in play.

## 1 **9 Supplementary Materials**

2 **Supplementary Materials by section** In Section 1 of the Supplementary Materials, we present  
3 a proof of Theorem 1. In Section 2, we discuss the asymptotic properties of our estimator,  
4 and present a theorem which provides an upper bound on the convergence rate. We prove  
5 this theorem in Section 3. In Section 4, we include tables illustrating the average practical  
6 performance for our approach (dtcode), umd, and scdensity, for the examples considered  
7 in the simulation study in the main paper. In Section 4.1, we discuss the effect of the  
8 number of basis elements on the final estimate. In Section 5, we include some examples  
9 of general shape-constrained density estimation beyond  $M$ -modality, like monotonicity, an  
10 upper bound on the number of modes, and so on. In Section 6, we include a simulation  
11 study for conditional density estimation. In Section 7, we discuss an application of shape-  
12 constrained density estimation to DNA methylation profiles.

## 13 **References**

- 14 Abramson, I. S. (1982). On bandwidth variation in kernel estimates—a square root law. *The annals*  
15 *of Statistics*, 1217–1223.
- 16 Bashtannyk, D. M. and R. J. Hyndman (2001). Bandwidth selection for kernel conditional density  
17 estimation. *Computational Statistics & Data Analysis* 36(3), 279–298.
- 18 Bickel, P. J. and J. Fan (1996). Some problems on the estimation of unimodal densities. *Stat.*  
19 *Sin.* 6(1), 23–45.



- 1 Birge, L. (1997). Estimation of unimodal densities without smoothness assumptions. *Ann.*  
2 *Stat.* 25(3), 970–981.
- 3 Brunner, L. J. and A. Y. Lo (1989). Bayes methods for a symmetric unimodal density and its  
4 mode. *Ann. Stat.* 17(4), 1550–1566.
- 5 Cheng, M.-Y., T. Gasser, and P. Hall (1999). Nonparametric density estimation under unimodality  
6 and monotonicity constraints. *J. Comput. Graph. Stat.* 8(1), 1–21.
- 7 Cordova, J., L. M. K. Sriram, A. Kocatepe, Y. Zhou, E. E. Ozguven, and R. Arghandeh (2018).  
8 Combined electricity and traffic short-term load forecasting using bundled causality engine.  
9 *IEEE Transactions on Intelligent Transportation Systems*, 1–11.
- 10 Dasgupta, S., D. Pati, and A. Srivastava (In press). A two-step geometric framework for density  
11 modeling. *Statistica Sinica*.
- 12 Efromovich, S. (2010). Orthogonal series density estimation. *Wiley Interdisciplinary Reviews:*  
13 *Computational Statistics* 2(4), 467–476.
- 14 Einbeck, J. and G. Tutz (2006). Modelling beyond regression functions: an application of mul-  
15 timodal regression to speed–flow data. *Journal of the Royal Statistical Society: Series C*  
16 *(Applied Statistics)* 55(4), 461–475.
- 17 Grenander, U. (1956). On the theory of mortality measurement: part ii. *Scand. Actuar. J.*
- 18 Hall, P. and L.-S. Huang (2002). Unimodal density estimation using kernel methods. *Statistica*  
19 *Sinica*, 965–990.

- 1 Harris, R. A., T. Wang, C. Coarfa, R. P. Nagarajan, C. Hong, S. L. Downey, B. E. Johnson,  
2 S. D. Fouse, A. Delaney, Y. Zhao, et al. (2010). Comparison of sequencing-based methods  
3 to profile dna methylation and identification of monoallelic epigenetic modifications. *Nature*  
4 *biotechnology* 28(10), 1097.
- 5 Izenman, A. J. (1991). Review papers: Recent developments in nonparametric density estimation.  
6 *J. Am. Stat. Assoc.* 86(413), 205–224.
- 7 Kwac, J., J. Flora, and R. Rajagopal (2014, Jan). Household energy consumption segmentation  
8 using hourly data. *IEEE Transactions on Smart Grid* 5(1), 420–430.
- 9 Lang, S. (2012). *Fundamentals of differential geometry*, Volume 191. Springer Science & Busi-  
10 ness Media.
- 11 Meyer, M. C. (2001). An alternative unimodal density estimator with a consistent estimate of the  
12 mode. *Stat. Sin.* 11(4), 1159–1174.
- 13 Pollard, K. S., M. J. Hubisz, K. R. Rosenbloom, and A. Siepel (2009). Detection of nonneutral  
14 substitution rates on mammalian phylogenies. *Genome research*.
- 15 Rao, B. L. S. P. (1969). Estimation of a unimodal density. *Sankhyā: The Indian Journal of*  
16 *Statistics, Series A (1961-2002)* 31(1), 23–36.
- 17 Srivastava, A. and E. P. Klassen (2016). *Functional and shape data analysis*. Springer.
- 18 Terrell, G. R. and D. W. Scott (1992). Variable kernel density estimation. *The Annals of Statistics*,  
19 1236–1265.

- 1 Turnbull, B. C. and S. K. Ghosh (2014). Unimodal density estimation using bernstein polynomi-  
2 als. *Comput. Stat. Data Anal.* 72, 13–29.
- 3 Ugray, Z., L. Lasdon, J. Plummer, F. Glover, J. Kelly, and R. Martí (2007). Scatter search and  
4 local nlp solvers: A multistart framework for global optimization. *INFORMS Journal on*  
5 *Computing* 19(3), 328–340.
- 6 Van Kerm, P. et al. (2003). Adaptive kernel density estimation. *Stata Journal* 3(2), 148–156.
- 7 Wahba, G. (1981). Data-based optimal smoothing of orthogonal series density estimates. *The*  
8 *annals of statistics*, 146–156.
- 9 Wasserman, L. and J. D. Lafferty (2006). Rodeo: Sparse nonparametric regression in high di-  
10 mensions. In *Advances in Neural Information Processing Systems*, pp. 707–714.
- 11 Wegman, E. J. (1970). Maximum likelihood estimation of a unimodal density, II. *Ann. Math.*  
12 *Stat.* 41(6), 2169–2174.
- 13 Wheeler, M., D. Dunson, and A. Herring (2017). Bayesian local extremum splines.  
14 *Biometrika* 104(4), 939–952.
- 15 Wolters, M. A. and W. J. Braun (2018). A practical implementation of weighted kernel density  
16 estimation for handling shape constraints. *Stat* 7(1), e202.