

---

# Multiresolution dictionary learning for conditional distributions

---

## Abstract

Nonparametric estimation of the conditional distribution of a response given high-dimensional features is a challenging problem. In many settings it is important to allow not only the mean but also the variance and shape of the response density to change flexibly with features, which are massive-dimensional with a distribution concentrated near a lower-dimensional subspace or manifold. We propose a multiresolution model based on a novel stick-breaking prior placed on the dictionary weights. The algorithm scales efficiently to massive numbers of features, and can be implemented efficiently with slice sampling. State of the art predictive performance is demonstrated for toy examples and a real data application.

Key words: Density regression; Dictionary learning; Manifold learning; Mixture of experts; Multiresolution stick-breaking; Nonparametric Bayes.

---

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

## 1. Introduction

Massive datasets are becoming a ubiquitous by-product of modern scientific and industrial applications. These data present statistical and computational challenges for machine learning because many previously developed approaches do not scale-up sufficiently. Specifically, challenges arise because of the ultrahigh-dimensionality, and relatively low sample size (the “large  $p$ , small  $n$ ” problem). Parsimonious models for such big data assume that the density in the ambient dimension concentrates around a lower-dimensional (possibly nonlinear) subspace. Indeed, a plethora of methodologies are emerging to estimate such lower-dimensional “manifolds” from high-dimensional data (??).

We are interested in using such lower-dimensional embeddings to obtain estimates of the conditional distribution of some target variable(s). This *conditional regression* setting arises in a number of important application areas, including neuroscience, genetics, and video processing. For example, one might desire automated estimation of a predictive density for a continuous neurologic *phenotype* of interest, such as intelligence or a creativity score, on the basis of available data for a patient including neuroimaging. The challenge is to estimate the probability density function of the phenotype *nonparametrically* based on an  $\mathcal{O}(10^6)$  dimensional image of the subject’s brain. It is crucial to avoid parametric assumptions on the density, such as Gaussianity, while allowing the density to change flexibly with predictors. Otherwise, one can obtain misleading predictions and poorly characterize predictive uncertainty.

There is a rich machine learning and statistical literature on conditional density estimation of a response  $y \in \mathcal{Y}$  given a set of features (predictors)  $x = (x_1, x_2, \dots, x_p) \in \mathcal{X}$ . Common approaches include hierarchical mixtures of experts (??), kernel methods (????), Bayesian finite mixture models (???) and Bayesian nonparametrics (?????).

In general, there has been limited consideration of scaling to large  $p$  settings, with the variational Bayes approach of (?) being a notable exception. For dimensionality reduction, Tran et al. follow a greedy variable

selection algorithm. Their approach does not scale to the sized applications we are interested in. For example, in a problem with  $p = 1,000$  and  $n = 500$ , they reported a CPU time of 51.7 minutes for a single analysis. We are interested in problems many orders of magnitude or more larger than this, and require a faster computing time while also accommodating flexible non-linear dimensionality reduction (variable selection is a limited sort of dimension reduction). To our knowledge, there are no nonparametric density regression competitors to our approach, which maintain a characterization of uncertainty in estimating the conditional densities; rather, all sufficiently scalable algorithms provide point predictions and/or rely on restrictive assumptions such as linearity.

In big data problems, scaling is often accomplished using divide-and-conquer techniques. Well known examples are classification and regression trees (CART) (?) and multivariate adaptive regression splines (MARS) (?). These algorithms fit surfaces to data by explicitly dividing the input space into a nested sequence of regions, and by fitting simple surfaces within these regions. Though these methods are appealing in providing a simple, flexible and interpretable mechanism of dimension reduction, it is well known that single tree estimates commonly have high variance and poor performance. There is a rich literature proposing improvements based on bagging (?), boosting (?) and random forests (?). Though these algorithms can substantially improve mean square error performance, computation can be expensive and performance degrades as dimensionality  $p$  increases.

In fact, a significant downside of divide-and-conquer algorithms is their poor scalability to high dimensional predictors. As the number of features increases, the problem of finding the best splitting attribute becomes intractable so that CART, MARS and multiple trees models cannot be efficiently applied. Also mixture of experts models become computationally demanding, since both mixture weights and dictionary densities are predictor dependent. In an attempt to make mixtures of experts more efficient, sparse extensions relying on different variable selection algorithms have been proposed (?). However, performing variable selection in high dimensions is effectively intractable: algorithms need to efficiently search for the best subsets of predictors to include in weight and mean functions within a mixture model, an NP-hard problem.

In order to efficiently deal with massive datasets, we propose a novel multiresolution approach which starts by learning a multiscale dictionary of densities, constructed as Gaussian within each set of a multiscale

partition tree for the features. This tree is efficiently learned in a first stage using a fast and scalable graph partitioning algorithm applied to the high-dimensional features (?). Expressing the conditional densities  $f(y|x)$  for each  $x \in \mathcal{X}$  as a convex combination of coarse to fine scale dictionary densities, the learning problem in the second stage is how to estimate the corresponding multiresolution probability tree. This is accomplished in a Bayesian manner using a novel multiresolution stick-breaking process, which allows the data to inform about the optimal bias-variance trade-off; weighting coarse scale dictionary densities more highly decreases variance while adding to bias if the finer scale structure is needed. This results in a model that allows borrowing information across different resolution levels and reaches a good compromise in terms of the bias-variance tradeoff. We show that the algorithm scales efficiently to massive numbers of features.

## 2. Setting

Let  $x \in \mathcal{X} \subseteq \mathbb{R}^p$  be a  $p$ -dimensional Euclidean vector-valued *predictor* random variable. Let  $f(x)$  denote the *marginal* probability density of  $x$ . We assume that  $f(x)$  concentrates around a lower-dimensional, possibly nonlinear, subspace  $\mathcal{M}$ . For example,  $\mathcal{M}$  could be a union of affine subspaces, or a smooth compact Riemannian manifold.

Let  $y \in \mathcal{Y} \subseteq \mathbb{R}$  be a real-valued *target* variable. We further assume that the *conditional* distribution is a function of only the position  $\mu$  of  $x$  within the subspace  $\mathcal{M}$ ,  $f(y|x) = f(y|\mu)$ . Let  $x$  and  $y$  be sampled from some true but unknown joint distribution. We would like to learn  $f(y|x)$ . We assume that we obtain  $n$  independently and identically sampled observations,  $(x_i, y_i)$ , for  $i \in \{1, 2, \dots, n\}$ . Our proposed model introduced in Section 3 is very general in accommodating an unknown density  $f(y|x)$  which changes nonparametrically according to the location of  $x$  in the lower-dimensional subspace. However, for exposition and testing of the model, it is useful to consider a simple example in which  $x$  lives on a smooth one-dimensional Riemannian submanifold embedded in  $\mathbb{R}^p$ , and  $y$  is a univariate Gaussian random variable whose mean and variance vary with the location of  $x$  along its geodesic.

We can formalize this simple example model as follows. Consider  $x_i \sim \mathcal{N}(\psi(\mu_i), \sigma^2 I_p)$ , where  $\Psi = \{\psi: \mathcal{M} \rightarrow \mathbb{R}^p\}$ ,  $\mu_i \in \mathcal{M}$ ,  $\sigma \in (0, \infty)$ ,  $I_p$  is the  $p \times p$  dimensional identity matrix, and  $\mathcal{N}(\cdot, \cdot)$  indicates a Gaussian distribution. Let  $\mathcal{M}$  be a smooth compact Riemannian manifold, such as the swissroll or the S-manifold. For simplicity, let us assume that  $\mathcal{M}$  is a curve. Let  $\psi(\mu) = 1_p \mu$  with  $1_p$  being a  $p$ -dimensional

vector with all elements equal to 1. Define the conditional  $f(y|x) = \mathcal{N}(\mu, g(\mu))$  for some positive function  $g(\cdot)$  defined on  $(0, \infty]$ . In other words, both the mean and standard deviation of  $y$  depend on the position of  $x$  along its geodesic. We will show in §5 that our construction facilitates a smooth estimate of the mean and variance of  $y$ , even though we are not explicitly smoothing, rather, the smoothness is induced via the model averaging over spatial resolution levels.

### 3. Model Specification

#### 3.1. Approach

Our approach proceeds in a two stage fashion as follows. We first learn a multiscale nonlinear partition tree of the feature data  $\{x_i\}_{i=1}^n$ , recursively partitioning  $\{x_i\}$  to obtain subsets that are increasingly homogeneous according to some metric. The coarsest scale contains all of the data, the next finer scale contains two or more clusters of observations having relatively similar  $x_i$  values, and so on until our convergence criteria are met (for example, the available sample size is exhausted so that few observations fall within each fine scale leaf partition set). Based on the multiscale partition, each value of  $x$  has an associated *path* through the tree encoding the set membership at each scale. Using the response data  $y_i$  for all subjects  $i$  in a given partition set, we estimate a dictionary density specific to that set and resolution level using Bayesian methods. The conditional density is then expressed as a convex combination of these multiresolution dictionary densities, with a posterior distribution on the weights learned under a multiresolution stick-breaking process.

#### 3.2. Model Structure

Suppose we define a multiscale partition of  $\mathcal{X}$ . Generation one corresponds to the entire  $\mathcal{X}$  denoted as  $\mathcal{X}^1$ . At generation two,  $\mathcal{X}^1$  is split into two or more mutually exclusive partition sets,  $\mathcal{X}^1 = (\mathcal{X}_1^2, \mathcal{X}_2^2, \dots)$ . Each subset is recursively partitioned into two subsets so that for a general partition level  $\ell$  the partition will be given by  $\mathcal{X}^\ell = (\mathcal{X}_1^\ell, \dots, \mathcal{X}_{\ell_S}^\ell)$ . Let us assume this process proceeds for  $L$  levels. Let  $(\ell, s)$  be the node associated to the  $s$ th subset at resolution level  $\ell$ . Let  $de(\ell, s)$  and  $an(\ell, s)$  be respectively the set of descendants and ancestors of node  $(\ell, s)$ . Let  $A_\ell(x) \in \{1, \dots, \ell_S\}$  be the subset at level  $\ell$  including predictor  $x$ , with  $A_1(x)$  equal to 1 by definition.

We characterize the conditional density  $f(y|x)$  as a convex combination of multiscale dictionary densities. At level one, the global parent density is denoted by  $f_1$ . For predictor value  $x$ , the dictionary density at

generation  $j$  is  $f_{B_j(x)}$  with  $B_j(x) = \{j, A_j(x)\}$ , for  $j = 2, \dots, k$ . Then,  $f(y|x)$  is defined as the convex combination of densities  $\{f_{B_j(x)}\}_{j=1}^k$  with weights  $\{\pi_{B_j(x)}\}_{j=1}^k$ , i.e.

$$f(y|x) = \sum_{j=1}^k \pi_{B_j(x)} f_{B_j(x)}(y), \quad (1)$$

where  $0 \leq \pi_{B_j(x)}$  and  $\sum_{j=1}^k \pi_{B_j(x)} = 1$ .

Each  $B(x)$  is a set encoding the path through the partition tree up to generation  $k$  specific to predictor value  $x$ . According to model (1), one observation can lie in subsets located at different resolution levels. This is critical in achieving a good compromise between bias and variance through borrowing information across different resolution levels. Though the proposed approach is reminiscent of a mixture of experts model (?), the two approaches are quite different, since under (1), neither mixture weights nor dictionary densities directly depend on predictors. This allows our model to scale efficiently to high dimensional predictors.

Now let us examine the implications of model (1). For two predictor values  $x$  and  $x'$  located close together, it is expected that the paths will be similar, which leads to similar weights on the dictionary densities. In the extreme case in which  $x$  and  $x'$  belong to the same leaf partition set, we have  $B(x) = B(x')$  and the path through the tree will be the same. Moreover, in this case, we will have  $f(y|x) = f(y|x')$  so that up to  $k$  levels of resolution the densities  $f(y|x)$  and  $f(y|x')$  are identical. If the paths through the tree differ only in the final generation or two, the weights will typically be similar but the resulting conditional densities will not be identical.

To derive mixture weights, a natural choice corresponds to a stick-breaking process (?). For each node  $B_j(x_i)$  in the binary partition tree, define a stick length  $V\{B_j(x_i)\} \sim \text{Beta}(1, \alpha)$ . The parameter  $\alpha$  encodes the complexity of the model, with  $\alpha = 0$  corresponding to the case in which  $f(y|x) = f(y)$ . We relate the weights in (1) to the stick-breaking random variables as follows:

$$\pi_{B_j(x)} = V\{B_j(x)\} \prod_{B_h \in an\{B_j\}} [1 - V\{B_h(x)\}],$$

with  $V\{B_k(x)\} = 1$  to ensure that  $\sum_{j=1}^k \pi_{B_j(x)} = 1$ . We refer to this prior as a *multiresolution stick-breaking process*.

#### 4. Estimation

The proposed approach is based on a two-stage algorithm where first the observations are allocated to different subsets in a tree fashion using an efficient partitioning algorithm and then, considering the partition as fixed, a multiresolution stick-breaking process is estimated. In practice, observations are partitioned applying metis (?), a fast multiscale technique used for graph partitioning. Though more complicated densities can be considered, dictionary densities  $f_{B_j}$  will be estimated by assuming a normal form, i.e.  $f_{B_j} = \mathcal{N}(\mu_{B_j}, \sigma_{B_j})$ . In particular, densities corresponding to a particular partition set will be estimated considering only observations belonging to that partition set. To be specific, for estimating density  $f_{B_j}(y)$ , we use the data  $\{y_i : x_i \in \mathcal{X}_{A_j}^j\}$ . We then conduct the analysis treating partition sets as fixed; this is critical for scalability to big  $p$ . On the surface, conditioning on a single tree seems overly restrictive, but using a second stage multiresolution probability model for the weights over the tree leads to inferences that are robust to the tree estimate; almost as if the tree itself were randomized.

Parameters involved in the dictionary densities can be estimated using either frequentist or Bayesian methods. Bayesian methods are appealing since they can avoid singularities associated with traditional maximum likelihood inference, the prior has an appealing role as a regularizer, and we can characterize uncertainty in dictionary learning through the resulting posterior. Hence, parameters involved in dictionary densities will be estimated through Bayesian methods and inference on stick breaking weights and dictionary density parameters will be carried out using the Gibbs sampler. For this purpose, introduce the latent variable  $S_i \in \{1, \dots, k\}$ , for  $i = 1, \dots, n$ , denoting the multiscale level used by the  $i$ th subject. Assuming data are normalized prior to analysis, we let  $\mu \sim \mathcal{N}(0, I)$  and  $\sigma = \mathcal{IG}(a, b)$  for the means and variances of the dictionary densities. Let  $n_{B_j}$  be the number of observations allocated to node  $B_j$ . Each Gibbs sampler iteration can be summarized in the following steps.

1. Update  $S_i$  by sampling from the multinomial full conditional with

$$\Pr(S_i = j | -) = \frac{\pi_{B_j(x_i)} f_{B_j(x_i)}(y_i)}{\sum_{h=1}^k \pi_{B_h(x_i)} f_{B_h(x_i)}(y_i)}$$

2. Update stick-breaking random variable  $V_{B_j(x_i)}$ , for  $j = 1, \dots, k$  and  $i = 1, \dots, n$ , from  $\text{Beta}(a_p, b_p)$  with  $a_p = 1 + n_{B_j}$  and  $b_p = \alpha + \sum_{B_h(x_i) \in de\{B_j(x_i)\}} n_{B_h(x_i)}$ .

3. Update  $(\mu_{B_j}, \sigma_{B_j})$  by sampling from

$$\mu_{B_j} \sim \mathcal{N}(\bar{y}_{B_j} n_{B_j} / \sigma_{B_j}, (1 + n_{B_j} / \sigma_{B_j})^{-1})$$

$$\sigma_{B_j} \sim \mathcal{IG} \left( a + n_{B_j} / 2, b + 0.5 \sum_{\{i: S_i = j\}} (y_i - \mu_{B_j(x_i)})^2 \right)$$

with  $\bar{y}_{B_j}$  being the average of the observation  $\{y_i\}$  allocated to node  $B_j$ .

#### 5. Simulation Studies

In order to assess the predictive performance of the proposed model, different simulation scenarios were considered. Let  $n$  be the number of observations,  $y \in \mathbb{R}$  the response variable and  $x \in \mathbb{R}^p$  a set of predictors. The Gibbs sampler was run considering 20,000 as the maximum number of iterations with a burn-in of 1,000. Gibbs sampler chains were stopped testing normality of normalized averages of functions of the Markov chain (?). Parameters  $(a, b)$  and  $\alpha$  involved in the prior density of parameters  $\sigma_{B_j}$ s and  $V_{B_j}$ s were set respectively equal to  $(3, 1)$  and 1.

First let us consider the simple toy example of §2. Figure 1(a) depicts the true mean and variance of  $y$  and our estimate as  $x$  moves along the geodesic. These estimates were obtained by performing leave-one-out prediction and considering the mean and variance of the predictive distribution of  $y_i$  as the mean and variance estimate of the  $i$ th observation. As the figure clearly shows our construction facilitates a smooth estimate of the mean and variance of  $y$ , even though we are not explicitly smoothing, rather, the smoothness is induced via the averaging over spatial scales.

In all other examples, predictors were assumed to belong to an  $r$ -dimensional space, either a lower dimensional plane or a non linear manifold, with  $r \ll p$ . For each synthetic dataset, the proposed model was compared with CART and lasso in terms of mean squared error. Competitors were limited by scalability considerations. In the first three simulation studies, the vector of predictors was assumed to lie close to a lower dimensional plane. In practice, predictors were modeled through a factor model, i.e.  $x_i = \Lambda \eta_i + \epsilon_i$  with  $\epsilon_i \sim \mathcal{N}_n(0, I_n)$ ,  $\Lambda$  being a  $(p \times r)$  matrix,  $\eta_i \sim \mathcal{N}_r(0, I_r)$  and  $r \ll p$ . The response  $y$  was assumed to be a function of the latent variable  $\eta$  so that the dependence between response and predictors was induced by the shared dependence on the latent factors. In all examples,  $\Lambda$  was assumed to be a sparse matrix with level of sparsity increasing with the number of columns and non zero elements of  $\Lambda$  drawn from a standard normal density. In the last two simulation studies, predictors

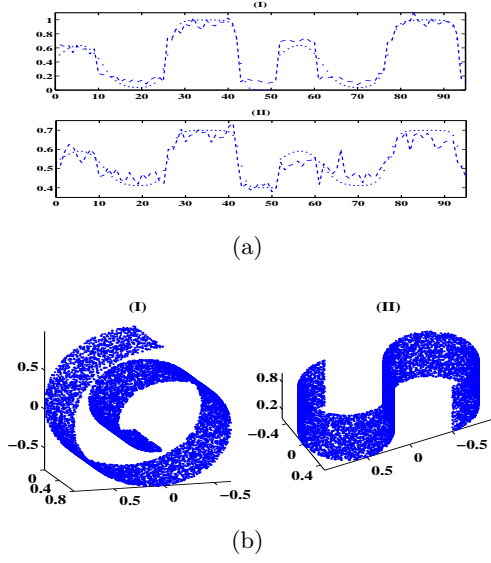


Figure 1. (a) Plot of mean and variance (I-II) for observations  $i = 1, \dots, 95$  (dot:true, dash:estimate); (b) Toy data examples: Swissroll (I) and S-Manifold (II) embedded in  $\mathbb{R}^3$

were assumed to lie close to the swissroll and the S-manifold (see figure 1(b)).

In the first simulation study, the response was assumed to be a linear function of predictors and they were jointly sampled from the above factor model. In the second simulation study, the response was drawn from a two components mixture of normals with mixture weights depending on the first latent factor, i.e.  $p = \exp\{\eta_1\}/(1 + \exp\{\eta_1\})$ , and components with location parameters  $(-2, 2)$  and unitary standard deviation. In the third simulation study, the response was drawn from a normal with mean and variance depending on the first latent factor as follows  $y \sim \mathcal{N}\{\eta_1^2 - \eta_1^3, \exp(1 - \eta_1)\}$ . In the last two simulation studies, predictors were drawn from the swissroll and the S-manifold, all two-dimensional manifolds but embedded in  $\mathbb{R}^{50}$ , while the response was sampled from a normal with mean equal to one of the coordinates of the manifold and standard deviation one.

Table 1 shows mean squared errors under the proposed approach, CART and lasso based on leave-one-out prediction. For each resolution level, the new observation was allocated to the set with closer center. As shown in table 1, CART performs worse than lasso only when the response is a linear function of predictors. However, in all data scenarios, our model is able to perform as well as or better than the model associated to the lowest mean squared error. Figure 2 shows the plot of CPU usage as a function of the number of

Table 1. Mean and standard deviations of squared errors under multiscale stick-breaking (MSB), CART and Lasso for sample size 50 and 100 for different simulation scenarios

SIM	$p$	$r$	MSB	CART	LASSO
(1)	1,000	5	1.09 (1.68)	2.29 (2.82)	1.09 (1.66)
(2)	10,000	5	0.55 (0.86)	0.55 (0.62)	0.99 (0.79)
(3)	5,000	5	0.78 (1.99)	0.83 (2.16)	0.84 (2.00)
(4)	50	2	0.80 (0.82)	1.00 (1.36)	1.01 (1.04)
(5)	50	2	0.60 (0.76)	0.64 (0.84)	1.01 (1.16)

Table 2. Real Data: Mean and standard deviations of squared error under multiscale stick-breaking (MSB), CART, Lasso and random forest (RF)

DATA	MODEL	MSE	$t_T$	$t_M$	$t_V$
(1)	MSB	0.56	100	1.1	0.02
	CART	1.10	87	0.9	0.01
	LASSO	0.63	200	2.8	0.17
	RF	0.57	7,817	78.2	0.59
(2)	MSB	0.76	690	20.98	2.31
	LASSO	1.02	5,836	96.18	9.66

features. This plot was obtained drawing  $(y_i, x_i)$ , for  $i = 1, \dots, 100$ , and  $x_i \in \mathbb{R}^p$  from the first simulation scenario considering different values of  $p$ . Clearly, our approach scales substantially better than competitors as the number of features increases.

Another important advantage of the proposed model is the possibility to obtain an estimate of the predictive density of the data. Figure 3 shows the estimated density of two data points sampled from the second simulation scenario. Clearly, the density function varies across different points and its estimate become closer to the true density as the number of observations in the training set increases.

## 6. Real Application

We assessed the predictive performance of the proposed method on two very different neuroimaging datasets. First, we consider a structural connectome dataset collected at the Mind Research Network. Data were collected as described previously (?). We inves-



550 tigated the extent to which we could predict creative  
 551 (as measured via the Composite Creativity Index (?)).  
 552 For each subject, we estimate a 70 vertex undirected  
 553 weighted brain-graph using the Magnetic Resonance  
 554 Connectome Automated Pipeline (?) from diffusion  
 555 tensor imaging data (?). We therefore let each  $x_i \in \mathbb{R}^p$   
 556 correspond to logarithm of each weighted edge; be-  
 557 cause our graphs are undirected and lack self-loops,  
 558 we have a total of  $\binom{70}{2} = 2415$  potential weighted  
 559 edges. The vector of covariates consists in the *JoVo*  
 560 *says: which log?* logarithm of the total number of  
 561 connections between all pairs of cortical regions, i.e.  
 562  $p = 2,415$ .

563 The second dataset comes from a resting-state func-  
 564 tional magnetic resonance experiment as part of the  
 565 Autism Brain Imaging Data Exchange. We selected  
 566 *JoVo says: which site* for analysis. Each brain-image  
 567 was processed using the Configurable Pipeline for  
 568 Analysis of Connectomes (?). For each subject we  
 569 computed a measure of normalized power at each voxel  
 570 called fALFF (?). fALFF is a highly nonlinear trans-  
 571 formation of the time-series data, previously demon-  
 572 strated to be a reliable property of such data. To en-  
 573 sure the existence of nonlinear signal relating these  
 574 predictors, we let  $y_i$  correspond to an estimate of over-  
 575 all head motion in the scanner, called mean framewise  
 576 displacement (FD) computed as described in (?). We  
 577 utilized a gray matter mask to consider only the voxels  
 578 with high probability of being gray matter. Thus, for  
 579 each of 56 subjects, we let  $x_i \in \mathbb{R}^{300,000}$  be the fALFF  
 580 of all gray matter voxels.

582 For the analysis, all variables were normalized by sub-  
 583 tracting the mean and dividing by the variance *JoVo*  
 584 *says: do you mean standard deviation? i imagine*  
 585 *so..* The same prior specification and Gibbs sampler  
 586 as in §5 was utilized. Table 2 shows mean and vari-  
 587 ance squared error based on leave-one-out predictions.  
 588 Variable  $t_T$  is the amount of time necessary to obtain  
 589 predictions for all subjects, while variables  $t_M$  and  $t_V$   
 590 are respectively the mean and the variance *JoVo says:*  
 591 *lets report standard deviation, it is on the same scale*  
 592 *as mean, unlike variance, which is related to mean*  
 593 *squared.* of amount of time necessary to obtain one  
 594 point predictions.

595 For the first data example, we compared our approach  
 596 (multiresolution stick-breaking; MSB) to CART, lasso  
 597 and random forests. For CART, lass, and random  
 598 forests, we use standard R packages downloaded from  
 599 CRAN *JoVo says: name the packages explicitly; ac-*  
 600 *tually put this up there when you first introduct them*  
 601 *for the simulations.* Table 2 shows that MSB out-  
 602 performs all the competitors in terms of mean square  
 603  
 604

error; this is in addition to yielding an estimate of the  
 entire conditional density for each  $y_i$ . It is also sig-  
 nificantly faster than random forests, the next closest  
 competitor, and faster than lasso. For this relatively  
 low-dimensional example, CART is reasonably fast.

For the second data application, given the huge di-  
 mensionality of the predictor space, we were unable to  
 get either CART or random forest to run to comple-  
 tion, yielding memory faults on our workstation *JoVo*  
*says: add computer details, how may cores, how much*  
*ram, etc..* We thus only compare performance to lasso.  
 As in the previous example, MSB outperforms lasso  
 in terms of predictive accuracy measured via mean-  
 squared error, and *significantly* outperforms lasso in  
 terms of computational time. Figure 4 shows the plot  
 of CPU time used to predict each one of the 56 sub-  
 jects involved in the experiment. The time needed to  
 compute quantities utilized in all subject predictions  
 was divided equally across subjects. Clearly, our ap-  
 proach is able to improve the computational time by  
 up to five orders of magnitude.

605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659

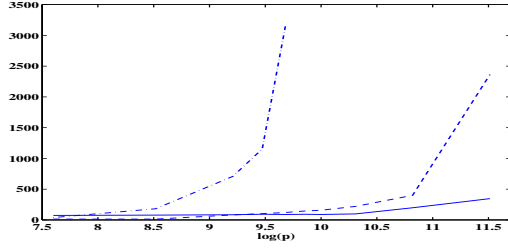


Figure 2. Elapsed CPU time (in seconds) for leave-one-out prediction based on 100 observations for MSB (solid), lasso (dash) and CART (dot-dash) for different number of predictors in log-scale

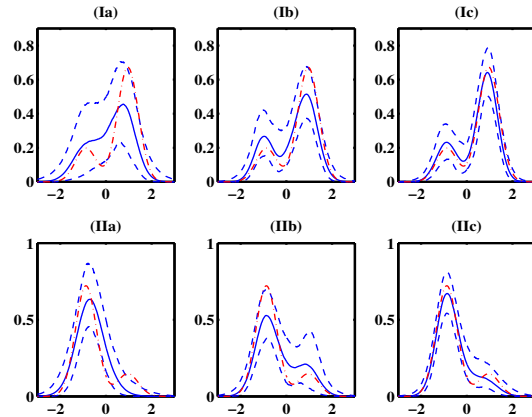


Figure 3. Plot of true (dashed-dotted line) and estimated (50th percentile: solid line, 2.5th and 97.5th percentiles: dashed lines) density for two data points (*I*, *II*) considering different training set size (a:50, b:100, c:150).

## Discussions

We have proposed a new model which should lead to substantially improved predictive and computational performance to learn the density of a target variable given a high dimensional vector of predictors. As shown the proposed two stage approach can scale substantially better than other existing algorithms to massive number of features. We have focused on Bayesian MCMC-based methods, but there are numerous interesting directions for ongoing research. Moreover, in addition to better predictive and computational performance, our methods easily extend to parallelized and distributed systems, which we will also explore in future work.

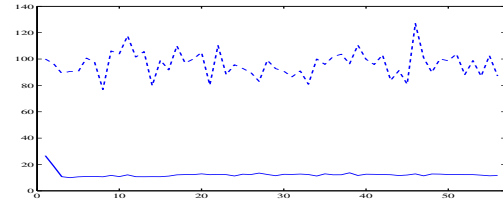


Figure 4. Plot of CPU time used to predict each one of the 56 measurement involved in experiment (2) under MSB (solid) and lasso (dash)