# Bandwidth selection for kernel based interval estimation of a density

Santanu Dutta[1*]

[1]*Tezpur University*

*Abstract:* It is always useful to have a confidence interval, along with a single estimate of the parameter of interest. We propose a new algorithm for kernel based interval estimation of a density, with an aim to minimize the coverage error. The bandwidth used in the estimator is chosen by minimizing a bootstrap estimate of the absolute value of the coverage error. The resulting confidence interval seems to perform well, in terms of coverage accuracy and length, especially for large sample size. We illustrate our methodology with data on the eruption durations for the Old Faithful geyser in USA. It seems to be the first bandwidth selector in the literature for kernel based interval estimation of a density.

*Key words*: Kernel based interval estimation, coverage error, bandwidth selection, bootstrap.
*AMS Subject Classification:* 62G07, 62G09, 62G20.

## 1. Introduction

We consider the problem of construction of a two sided confidence interval for $f(x_0)$, where $f$ is the unknown density generating the given data and $x_0$ is a given design point. A density function may be arbitrarily specified at a point $x_0$. This technical difficulty is overcome by assuming that $f$ is a continuous function. In the sequel we assume that $f$ is continuous and $x_0$ is an interior point of the support of $f$.

One of the most well known estimators of $f$ is a kernel density estimator (KDE) defined as follows.

Let $X_1,....,X_n$, be independent and identically distributed random variables with an unknown density $f(\cdot)$. The kernel density estimator of $f$ based on the kernel $K(\cdot)$ and bandwidth $h \equiv h_n$, is defined as

$$\hat{f}_n(y) \equiv f(y,h) = \frac{1}{nh}\sum_{i=1}^{n} K\left(\frac{y-X_i}{h}\right) \tag{1.1}$$

where $h \to 0$ and n$h \to \infty$ as n $\to \infty$. The problem of data based selection of $h$ for estimating

$f(x_0)$ using $\hat{f}_n(x_0)$ has been well studied. See for instance, [1] and [3] among most recent.

---

* Corresponding author.

In contrast, far less seems to be known regarding the choice of $h$ for constructing a confidence interval for $f(x_0)$ using $\hat{f}_n(x_0)$. For instance, in [1] the authors mention that there seems to be no automatic method for practical interval estimation for $f(x_0)$ available in the literature. From the simulation study in [6] we see that the bandwidth which is appropriate (in terms of coverage accuracy) for confidence interval construction is not easy to determine. No data based method for selecting such a $h$ was suggested by the author. Chen proposed empirical likelihood confidence intervals for density estimation, but again no bandwidth selection method was provided (see [2]). Fiorio has discussed two programs, viz. "asciker" and "bsciker" in Stata, to compute asymptotic and bootstrap confidence intervals for kernel density estimation. However these programs assume that the search for the correct bandwidth has been performed beforehand (see page 173 in [8]). Therefore these algorithms cannot be used to determine the appropriate amount of smoothing for kernel based interval estimation. In this paper we propose an algorithm for data based choice of $h$ with an aim to minimize the coverage error of the resulting confidence interval.

A kernel based confidence interval for $f(x_0)$ crucially depends on the approximations of the quantiles of the sampling distribution of $S = \hat{f}_n(x_0) - E[\hat{f}_n(x_0)]/\sigma^2$ and the bias $b = E[\hat{f}_n(x_0)] - f(x_0)$, where $\sigma^2$ is an estimated standard deviation of $\hat{f}_n(x_0)$.

We use the following $\sigma^2 \equiv \hat{\sigma}^2(h)$ proposed by Hall in [6].

$$\hat{\sigma}^2(h) = \frac{1}{nh}\left[\frac{1}{nh}\sum_{i=1}^{n} K\left(\frac{x_0 - X_i}{h}\right) - h\hat{f}_n(x_0)^2\right] \qquad (1.2)$$

The bias $b$ is not negligible even for a bandwidth minimizing the mean squared error. There are two approaches to tackle the bias $b$, viz. either to estimate the bias explicitly, or to reduce it substantially by under smoothing (see [6]). In [6], the author showed that under smoothing method produces confidence intervals with greater coverage accuracy than those obtained by explicit bias correction. There are several other practical advantages of the under smoothing method. For instance, in the under smoothing approach no estimator of the bias is required.

In the under smoothing approach we essentially construct a confidence interval for $E(\hat{f}_n(x_0))$ using a small value of the bandwidth, such that the same interval can be used to perform inference on $f(x_0)$ Horowitz suggested to perform under smoothing with $n^{-k}$, $k>1/(2r+1)$, where r is the kernel order (see [7]). From [6] we see that the bandwidth which minimizes the asymptotic coverage error of a two sided under smoothed interval is of the form $h = H_{n^{-1/(r+1)}}$, where $H$ is a constant depending on $f^{(r)}$. However, Hall pointed out that substantial under smoothing is not a practical proposition. He suggested to use $h = c 1.05 \backslash \hat{\gamma} n^{-1/5}$ for under smoothing $\hat{f}_n(x_0)$, where $0 < c < 1$. $\hat{\gamma}$ is the sample standard deviation. The values of $c$ which give good coverage accuracy for given $x_0$, $n$ and distribution are not easy to determine empirically (see numerical results in [6]). We provide a solution to this problem.

Given $X_1, ...., X_n$ and a bandwidth h, a two sided under-smoothed bootstrap $1-\alpha$ confidence interval of $f(x_0)$ is defined as

$$I(1-\alpha) \equiv I(1-\alpha)(h) = (\hat{f}_n(x_0,\ h) - \hat{\sigma}(h)\hat{u}_{1-\alpha/2},\ \hat{f}_n(x_0,\ h) - \hat{\sigma}(h)\hat{u}_{\alpha/2}),$$

where $\hat{u}_\alpha$ is the $\alpha$th quantile of a bootstrap approximation of the sampling distribution of $S$. It is in fact the $\alpha$th quantile of the conditional distribution of $S^* = (\hat{f}^*{}_n(x_0, h) - \hat{f}_n(x_0, h))/\hat{\sigma}^*$, given $X_1, ...., X_n$. $\hat{f}^*{}_n(x_0)$, $\hat{\sigma}^*$ are the versions of $\hat{f}_n(x_0)$, and $\hat{\sigma}$, obtained by replacing $X_1, ...., X_n$ by Efron's (1979) classical bootstrap re-sample $X_1^*, ...., X_n^*$. in (1.1) and (1.2).

The (exact) coverage probability of $I(1 - \alpha)$ is defined as

$$\beta(1 - \alpha) \equiv \beta(1 - \alpha)(h) = P(f(x_0) \in I(1 - \alpha)).$$

[6] suggested to select h with an aim to minimize the absolute value of the coverage error, viz. $CE = |\beta(1 - \alpha) - (1 - \alpha)|$. However $\beta(1 - \alpha)$ is a function of the unknown $f$. So for practical data based choice of h, CE has to be estimated based on $X_1, ...., X_n$. Using Efron's (1979) classical bootstrap method we propose an estimate of the $CE$ and it is minimized (with respect to $h$) for data based choice of the bandwidth. Let $\hat{h}$ denote the proposed data based bandwidth. The details of our proposal are given in Section 2.

The exact coverage probability $\beta(1 - \alpha)(\hat{h})$, of the confidence interval using $\hat{h}$, is hard to compute. However for any given $f$, we can approximate the coverage probability using Monte-Carlo simulations. In a simulation study, in Section 3, we compute the Monte-Carlo estimates of $\beta(1 - \alpha)(\hat{h})$ for different choices of $f$, $x_0$ and $n$. We also report the average width and the variance of the widths of the confidence intervals. These results are compared with the findings in [6].

## 2. Our proposal

Given $X_1, ...., X_n$ and $h$, we propose a bootstrap estimate $\beta^*(1 - \alpha)$ of the coverage probability $\beta(1 - \alpha)$ as follows

$$\beta^*(1 - \alpha) \equiv \beta^*(1 - \alpha)(h) = P^*(\hat{f}_n(x_0, h) \in I^*(1 - \alpha)),$$

where

$$I^*(1 - \alpha) = (\hat{f}_n^*(x_0, h)| - \hat{\sigma}^*(h)\hat{u}^*_{1-\alpha/2}, \ \hat{f}_n^*(x_0, h) - \hat{\sigma}^*(h)\hat{u}^*_{\alpha/2}).$$

Given $X_1, ...., X_n$, let $X_1^*, ...., X_n^*$ be a simple random sample drawn with replacement (srswr) from the empirical distribution. As mentioned earlier $\hat{f}^*{}_n(x_0)$, $\hat{\sigma}^*$ are the bootstrap versions of $\hat{f}_n(x_0)$ and $\hat{\sigma}$. $P^*$ denotes the conditional probability, given $X_1, ...., X_n$. $\hat{u}^*_\alpha$ is a bootstrap version of the statistic $\hat{u}_\alpha$. In fact $\hat{u}^*_\alpha$ is the $\alpha$th quantile of the conditional distribution of $S^{**} = (\hat{f}^{**}{}_n(x_0, h) - \hat{f}^*{}_n(x_0, h)) / \hat{\sigma}^{**}$, given $X_1^*, ...., X_n^*$ and $h$. $\hat{f}^{**}{}_n(x_0, h)$ and $\hat{\sigma}^{**}$ are obtained by replacing $X_1, ...., X_n$ in (1.1) and (1.2) by $X_1^{**}, ...., X_n^{**}$ which is a second stage re-sample drawn with replacement from $X_1^*, ...., X_n^*$.

$\beta^*(1 - \alpha)$ is a function of the bandwidth $h$. We define a bootstrap estimator $CE$ of the coverage error as follows

$$\widehat{CE} \equiv \widehat{CE}(h) = |\beta^*(1 - \alpha)(h) - (1 - \alpha)|.$$

We minimize $\widehat{CE}$ with respect to $h$ for data based bandwidth selection. The resulting random $\hat{h}$ is defined as follows

$$\hat{h} = \mathrm{argmin}_{h \in J_n} \widehat{CE}(h), \tag{2.1}$$

where $J_n$ is a compact interval with endpoints equal to scale invariant bandwidths, which are smaller than the bandwidth minimizing the MISE. As mentioned earlier, Hall suggested to use $\hat{h} = 1.05c\hat{\gamma}\, n^{-1/5}, 0< c \leq 1$, for under smoothing (see [6])}.  Motivated by this proposal we use

$$J_n = \left[c_1 1.05\hat{\gamma}\, n^{-1/5}, \ \ c_2 1.05\hat{\gamma}\, n^{-1/5}\right], 0 < c_1 < c_2 \leq 1$$

Hall considered a wide range of values of $c$ varying from 0.1 to 1, and showed that widely different values of $c$ are appropriate under different circumstances (see Table 1 in page 687 in [6]). Motivated by this, we use $c_1$=0.1 and $c_2$=1. With these choices of $c_1, c_2$ , $J_n$ covers all the under smoothing bandwidths considered by Hall in the simulation study in [6].

The proposed two sided under smoothed bootstrap $1- \alpha$ confidence interval of $f(x_0)$ is defined as

$$I(1 - \alpha)(\hat{h}) = (\hat{f}_n(x_0, \ \hat{h}) - \hat{\sigma}(\hat{h})\hat{u}_{1-\alpha/2}, \ \hat{f}_n(x_0, \ \hat{h}) - \hat{\sigma}(\hat{h})\hat{u}_{\alpha/2}). \tag{2.2}$$

## 2.1  Some computational details

### 2.1.1 Computation of  $\hat{u}_\alpha^*$

Given $X_1,...., X_n$ and $h$, we compute $\hat{u}_\alpha$ as follows.

We draw $B_2$ bootstrap re-samples. For each re-sample we compute $S^*$. There are $B_2$ values of  $S^*$ corresponding to the re-samples. Now $\hat{u}_\alpha$ is the $\alpha$th sample quantile based on these $B_2$ values.

### 2.1.2 Computation of $\hat{u}_\alpha^*$

Let $X_1^*,....,X_n^*$ be a bootstrap re sample drawn from  $X_1,...., X_n$.  Based on $X_1^*,....,X_n^*$, we compute $\hat{u}_\alpha^*$ as follows.

We generate $B_2$ second stage re-samples from $X_1^*, ....,X_n^*$, and compute the values of $S^{**}$ based on the $B_2$ second stage re-samples. The $\alpha$th sample quantile of these $B_2$ values of  $S^{**}$  is a Monte Carlo approximation to $\hat{u}_\alpha^{**}$.

### 2.1.3 Computation of $\beta^*(1- \alpha)$ $(h)$

Given $X_1,...., X_n$ and $h$ the computation of $\beta^*(1- \alpha)$ (h) involves the following  steps.

(i)    Generate $B_1$ re-samples, each of size $n$, by simple random sampling with replacement (srswr) from $X_1,...., X_n$, and compute $\hat{f}^*_{\ n}(x_0, h)$, $\sigma^*(h)$ for each re-sample.

(ii) From each re-sample, we further generate $B_2$ second stage re-samples by srswr. Using these second stage re-samples we compute $\hat{u}^*_{\alpha/2}$ and $\hat{u}^*_{1-\alpha/2}$ by the procedure mentioned above.

(iii) Using $\hat{f}^*_n(x_0, h)$, $\sigma^*(h)$, $\hat{u}^*_{\alpha/2}$ and $\hat{u}^*_{1-\alpha/2}$, we compute $I^*(1-\alpha)$ for each (1st stage) re-sample. There are $B_1$ such intervals corresponding to the $B_1$ first stage re-samples.

(iv) The Monte-Carlo estimate of $\beta^*(1-\alpha)$ $(h)$ is equal to the number of the intervals (obtained in step (iii)) containing $\hat{f}_n(x_0, h)$ divided by $B_1$.

## Remark 1.

1. As mentioned earlier $I^*(1-\alpha)$ is a two sided confidence interval for $E\left(\hat{f}_n(x_0)\right)$. The above mentioned algorithm essentially imitates the Mone-Carlo (MC) method of approximating the exact coverage probability of $\beta(1-\alpha)(h)$, for any given $f$ and $h$. In the MC method we draw random samples from a given distribution, and for each sample we compute $I(1-\alpha)$ by the re-sampling method described earlier. The MC estimate of $\beta(1-\alpha)(h)$ is the number of the intervals containing $E\left(\hat{f}_n(x_0)\right)$ divided by the number of random samples drawn. We imitate this procedure, replacing the actual distribution by the empirical distribution.

We note that $\hat{f}_n(x_0) = E^*\left(\hat{f}^*_n(x_0)\right)$, where $E^*$ denotes the expectation with respect to the empirical distribution. So the bootstrap version of $I(1-\alpha)$ is a confidence interval for $\hat{f}_n(x_0)$, given $X_1, ...., X_n$. In our method the 1st stage re-samples, drawn from the empirical distribution, mimic the role played by the random samples drawn from the actual distribution in the MC method.

2. We use the same 1st stage re-samples and 2nd stage re-samples (obtained by re-sampling each 1st stage re-sample in step [ii] of the above algorithm) to compute $\beta^*(1-\alpha)$ $(h)$ for different values of $h$, as required in a numerical minimization algorithm. This feature reduces the computational burden.

3. Given a confidence interval, Monte-Carlo approximation of its coverage probability essentially involves estimating an average of a random function using Monte-Carlo simulations. From [5] we see that much larger number of Monte-Carlo re-samples are required for approximating a bootstrap quantile estimator accurately, than the same required for approximating a bootstrap estimator of the expectation of some random function. Therefore we use different number of re-samples, viz. $B_2$ and $B_1$, to approximate the bootstrap estimators of the quantiles and the coverage probability by Monte-Carlo method.

## 2.2 Monte Carlo sample size for bootstrap-resampling

From [10] we see that the selection of appropriate $B_1$ and $B_2$ are not easy problems. As a rule of thumb, [5] suggested that for Monete-Carlo approximation of bootstrap moment

estimators the number of bootstrap re-samples should be 50 to 200.For approximating bootstrap quantile estimators the number of bootstrap re-samples should be at least 1000 (see [5] ) . We use this rule of thumb, and use $B_1$=200, $B_2$=1000.

## 3.   Simulation

Hall conducted simulations to study the effect of the choice of h on the cover-age probability of an under smoothed bootstrap confidence interval $I(1 - \alpha)(h)$was examined for six combinations of $f$ and $x_0$ (see [6]). The author used $h = c1.05\hat{\gamma}n^{-1/5}$, where $0 < c \leq 1$, for under smoothing the density estimator. In his simulations $f$ equals to the $N(0, 1)$ density and the $(1/2)N(0, 1)+(1/2)N(3, 1)$ density, and $x_0$ equal to 0, 0.75 and 1.5. The notation $pN(\mu_1, \sigma_1^2) + (1 - p)N(\mu_2,\sigma_2^2)$ represents a two component mixed normal distribution, where $\mu_i$, $\sigma_i^2$ are the mean and variance of the ith mixing component. For both these test den-sities, $x_0 = 0$ is the peak of the density. Hall reported the Monte Carlo estimates of the exact coverage probability $\beta(1 - \alpha)(h)$, along with the average and stan-dard deviation of the interval length. It was observed that the coverage accuracy of the confidence interval for f at the peak was less than the same at other point.

In [1], the authors considered the problem of interval estimation of $f(0)$, where $f$ is a standard normal density. From their simulations (page 513, in [1] we see that neither the coverage error nor the length of their 95 percent interval seem to decrease as $n$ is increased more than two times. This is perhaps due to the fact that random bandwidth proposed by Chan Lee and Peng is suitable for point estimation of $f$ at $x_0$ . In [7], the author pointed out that nonparametric point estimation and interval estimation are different tasks that require different degrees of smoothing.

In this section we study effect of the proposed random bandwidth $\hat{h}$ on the coverage probability and the average length of $I(1 - \alpha)$, for different choices of $f$ and $x_0$ and α =0.05. We consider the above mentioned choices of $f$ and $x_0$ as in [6]. Both these densities are unimodal, with peak at $x_0$= 0. In addition we consider two more test densities, viz. $f$ equal to the $(1/2)N(-1, 1/2)+(1/2)N(1, 1/2)$ density and the gamma(2,1) density. For the $(1/2)N(-1, 1/2) + (1/2)N(1, 1/2)$ density there are two peaks of same height at −1 and 1, and a trough at 0. We estimate this density at $x_0$ equal to 0 and 1. For the gamma density peak occurs at 1. We estimate the height of the gamma density and $x_0$ equal to 1 and 4.474, which is the 95th percentile. To compute the Monte-Carlo estimate of the coverage probability of a confidence interval we draw m random samples of a specific size from a test distribution, and compute the confidence interval for each sample. So there are $m$ such intervals. The Monte-Carlo estimate of the coverage probability is equal to number of intervals containing $f(0)$, divided by m. In Table 1 we use $c_1$ =0.1 and $c_2 = 1$.

In Table 2 we report the Monte-Carlo estimates of the coverage probability, average length and variance of the confidence intervals using $h = c1.05\hat{\gamma}n^{-1/5}$ , for different choices of $c$ and $f$ equal to the $(1/2)N(-1, 1/2) + (1/2)N(1, 1/2)$ density and the gamma(2,1) density. If the mean or the variance of the length of the confidence interval exceeds 100, we write "large".

In Table 1 we report the Monte-Carlo estimate of the coverage probability, average length and variance of the proposed confidence interval $I(1-\alpha)(\hat{h})$, in (2.2), for 10 combinations of $f$ and $x_0$. We compute each estimate for n = 50 and n = 100. To compute Monte-Carlo estimate we draw m = 300 samples from each test density. We have the following observations.

Table 1: Monte Carlo estimates of $\beta(1-\alpha)(\hat{h})$ for $h$ eaual to $\hat{h}$ and $\alpha = 0.05$

| *Density* | $(x_0, n)$ | Coverage Probability | Interval Width average (variance) |
|---|---|---|---|
| N(0,1) | (0, 50) | 0.90 | 0.371 (0.014) |
|  | (0, 100) | 0.96 | 0.151 (0.002) |
|  | (0.75, 50) | 0.91 | 0.381 (0.013) |
|  | (0.75, 100) | 0.958 | 0.239 (0.006) |
|  | (1.5, 50) | 0.88 | 0.221 (0.007) |
|  | (1.5, 100) | 0.935 | 0.143 (0.002) |
| (1/2)N(-1, 1/2) + (1/2)N(1, 1/2) | (0, 50) | 0.90 | 0.229 (0.009) |
|  | (0, 100) | 0.91 | 0.167 (0.003) |
|  | (1, 50) | 0.90 | 0.384 (0.033) |
|  | (1, 100) | 0.91 | 0.295 (0.005) |
| (1/2)N(0, 1) + (1/2)N(3, 1) | (0, 50) | 0.924 | 0.179 (0.003) |
|  | (0, 100) | 0.935 | 0.129 (0.001) |
|  | (0.75, 50) | 0.97 | 0.162 (0.002) |
|  | (0.75, 100) | 0.962 | 0.117 (0.001) |
|  | (1.5, 50) | 0.915 | 0.160 (0.012) |
|  | (1.5, 100) | 0.94 | 0.112 (0.001) |
| gamma(2,1) | (1, 50) | 0.87 | 0.306 (0.011) |
|  | (1, 100) | 0.965 | 0.255 (0.004) |
|  | (4.474,50) | 0.84 | 0.081 (0.001) |
|  | (4.474,100) | 0.88 | 0.071 ( 0.002) |

(i)     The confidence interval $I(1-\alpha)(\hat{h})$, using the proposed random bandwidth $\hat{h}$ in (2.1), seems to perform consistently. The coverage error, the mean and the variance of the interval length seem to reduce as sample size is increased for all choices of $f$ and $x_0$.

(ii)    From the simulation study in [6] and our Table 2, we see that the coverage probability and length of the confidence intervals using $h = c1.05\hat{\gamma}n^{-1/5},\backslash\ 0 < c \le 1$, can vary widely depending on estimation point $x_0$ and $c$.

(iii)   In contrast, the simulations in Table 1 indicate that  for a given distribution the coverage accuracy of the   confidence interval using  $\hat{h}$ does not seem to vary drastically with the change in  $x_0$, especially for n=100. This is due to the fact that proposed bandwidth selector is a function of the estimation point $x_0$, and so the resulting bandwidth  $\hat{h}$ automatically adjusts the amount of smoothing depending on $x_0$.

(iv) From the simulations in [6] we see that for $f$ equal to the $(1/2)N(0, 1) + (1/2)N(3, 1)$ density and $x_0$ equal to the peak, the coverage probability of the under smoothed confidence interval is poor especially for $c > 0.5$ in $h = c1.05\hat{\gamma}n^{-1/5}$. From our Table 2 we see that a similar observation is also true for $x_0$ equal the trough between the two peaks of the $(1/2)N(-1, 1/2) + (1/2)N(1, 1/2)$ density. Hall pointed out that the coverage error of confidence interval for estimation $f$ at the peak is in general higher than the same at other points, as the bias in a kernel density estimator is more pronounced at a peak. We observe that the same argument is also true for $x_0$ equal to a trough. Moreover from Table 2 we see that while estimating the gamma density at the peak the under smoothed confidence interval using  $h = c1.05\hat{\gamma}n^{-1/5}$ performs poorly for every choice $c$.

However, simulations in Table 1 suggest that the proposed confidence interval $I(1 - \alpha)(\hat{h})$performs  well in estimating $f$ at  the peak as well as the trough, in terms of the coverage accuracy, especially for n=100 and irrespective of $f$.

(v) From the simulations in [6] and our Tables 1 and 2, we see that the mean and the variance of the length proposed confidence interval compares well with the lengths of the corresponding confidence intervals using $h = c1.05\hat{\gamma}n^{-1/5}$ in [6].

## 3.1 Faithful data analysis

A well known data set in the context of density estimation is the data on the durations (in minutes) of eruptions for the Old Faithful geyser in Yellowstone National Park, Wyoming, USA. It is available in the R software (see data set ``faithful'' in R). The histogram based on the frequency density of the raw data is plotted in Figure 1.  We construct the 95 percent confidence intervals at 30 equidistant grid points using the proposed method. These upper and lower limits of the confidence intervals are marked as red and blue ``bubbles'' in Figure 1. We also plot the kernel density estimates using the plug-in bandwidths proposed by Sheather and Jones (see [9]) and the least square cross validation bandwidth. The cross validation and plug-in density estimates are numbered as 1 and 2 in Figure 1.  We observe the following.

The data is strongly bi-modal. The upper limits of the 95 percent confidence intervals seem to close to the frequency density of the raw data at the grid points (see the red "bubbles" in Figure 1"). The left-peak in the cross validation based curve is taller than the same in the plug-in curve.  Both the density estimates are within the 95 percent confidence interval near the two peaks. It seems that the left peak of the underlying density can be taller than the same in the plug-in based curve. The cross validation based density estimate seems to be reasonable near the left-peak. The cross validation curve is always within the confidence intervals at the grid points. The plug-in curve seems to lie outside confidence interval at the grid points in the left tail region. So for this data set the cross validation density estimate seems to be a more reasonable fit.

Clearly the confidence intervals, along with the point estimates of the density, enable a more detailed analysis of the data than that based on a single density estimate.
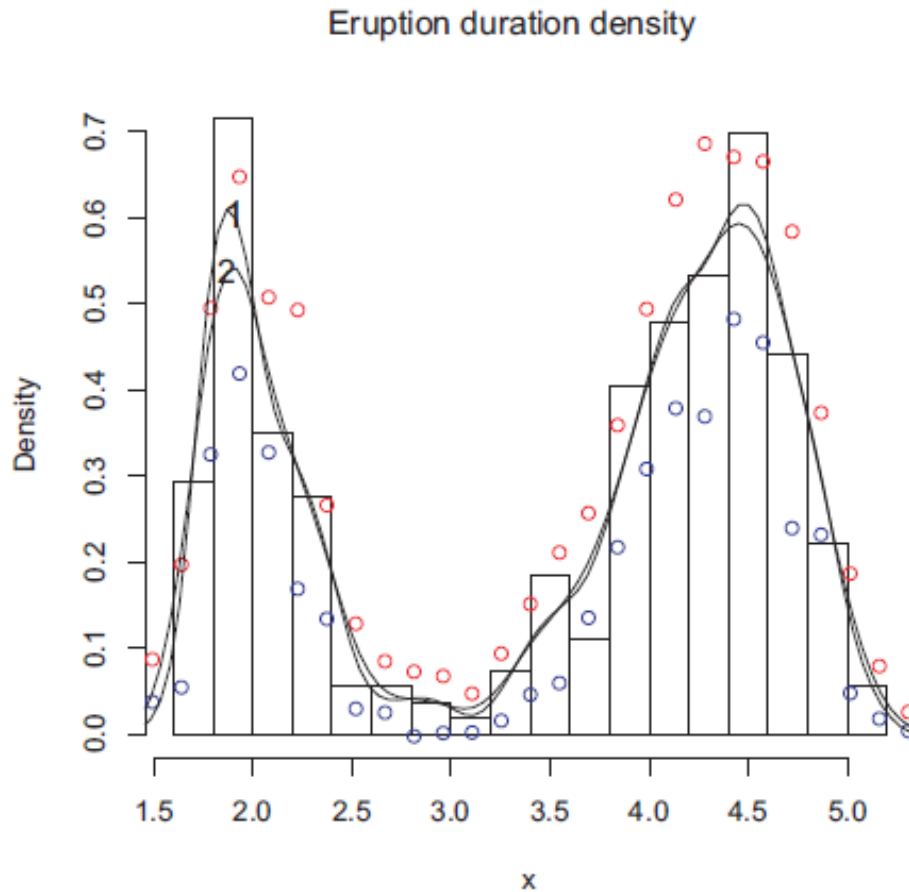


Figure 1: Fig 1: the 95 percent confidence bands and the density estimates using plug-in bandwidth (the red curve) and the least square cross validation bandwidth (the blue curve) based on the eruption durations data.

## Final Remarks

From the above simulation study it appears that the confidence interval $I(1 - \alpha)(\hat{h})$ in (2.2) performs well for all the test densities, especially for n=100. Simulations in our Table 2 suggest that if $f$ is a density with positive support and $x_0$ is the peak, the under smoothed confidence interval for $f(x_0)$ using $h = c1.05\hat{\gamma}n^{-1/5}$ performs poorly for all the different choices of $c$ mentioned in [6]. In contrast, the coverage error or the average length of $I(1 - \alpha)(\hat{h})$ does not seem to vary drastically for different choices of $x_0$. So the proposed bandwidth selector can be recommended safely for interval estimation of $f(x_0)$, especially for large sample size. We observe that using a confidence band in conjunction with the usual global density estimates

more detailed information can be extracted from the faithful eruption duration data, than that obtained using a single density estimate.

Table 2: Monte Carlo estimates of $\beta(1-\alpha)(\hat{h})$ for $h = c1.05\hat{\gamma}n^{-1/5}$ for different values $c$

| Density | $c$ | $(x_0, n)$ | Coverage Probability | Interval width mean (variance) |
|---|---|---|---|---|
| (1/2)N(-1, 1/2) + (1/2)N(1, 1/2) | 0.1 | (0, 50) | 0.95 | 2435.31(large) |
| | | (0, 100) | 0.975 | 15.744 (large) |
| | | (1, 50) | 0.975 | 86.57 (large) |
| | | (1, 100) | 0.97 | 0.850 (0.219) |
| | 0.2 | (0, 50) | 0.965 | 243 (large) |
| | | (0, 100) | 0.96 | 1.672 (19.850) |
| | | (1, 50) | 0.98 | 0.616 (0.022) |
| | | (1, 100) | 0.96 | 0.432 (0.004) |
| | 0.3 | (0, 50) | 0.95 | 0.879 (1.656) |
| | | (0, 100) | 0.965 | 0.272 (0.019) |
| | | (1, 50) | 0.955 | 0.455 (0.0038) |
| | | (1, 100) | 0.975 | 0.323 (0.0013) |
| | 0.5 | (0, 50) | 0.865 | 0.233 (0.0026) |
| | | (0, 100) | 0.87 | 0.156 (0.0006) |
| | | (1, 50) | 0.965 | 0.303 (0.001) |
| | | (1, 100) | 0.9 | 0.227 (0.0005) |
| | 0.75 | (0, 50) | 0.415 | 0.161 (0.0005) |
| | | (0, 100) | 0.32 | 0.120 (0.0002) |
| | | (1, 50) | 0.755 | 0.220 (0.0003) |
| | | (1, 100) | 0.62 | 0.167 (0.0002) |
| | 1 | (0, 50) | 0.01 | 0.125 (0.0003) |
| | | (0, 100) | 0.01 | 0.097 (0.0001) |
| | | (1, 50) | 0.285 | 0.172 (0.0003) |
| | | (1, 100) | 0.225 | 0.132 (0.0001) |
| Gamma(2,1) | 0.1 | (0, 50) | 0.83 | large(large) |
| | | (0, 100) | 0.75 | " |
| | | (4.474, 50) | 0.975 | " |
| | | (4.474, 100) | 0.955 | " |
| | 0.2 | (0, 50) | 0.425 | large(large) |
| | | (0, 100) | 0.195 | " |
| | | (4.474, 50) | 0.975 | " |
| | | (4.474, 100) | 0.965 | " |
| | 0.3 | (0, 50) | 0.125 | 4.019(large) |
| | | (0, 100) | 0.01 | 0.198 (0.024) |
| | | (4.474, 50) | 0.97 | large (large) |
| | | (4.474, 100) | 0.965 | " |

| | 0.5 | (0, 50) | 0.01 | 0.199 (0.024) |
|---|---|---|---|---|
| | | (0, 100) | 0.01 | 0.103 (0.0004) |
| | | (4.474, 50) | 0.96 | large (large) |
| | | (4.474, 100) | 0.975 | " |
| | 0.75 | (0, 50) | 0.0 | 0.123 (0.006) |
| | | (0, 100) | 0.01 | 0.088 (0.0001) |
| | | (4.474, 50) | 0.98 | 194.5 (large) |
| | | (4.474, 100) | 0.955 | 0.096 (0.006) |

## Reference

[1] Chan, N.H., Lee, T.C.M., and Peng, L. (2010). On nonparametric local inference for density estimation. *Comuptational Statistics and Data Analysis* **54**: 509-515.

[2] Chen, S.X., 1996. Empirical likelihood confidence intervals for nonparametric density estimation. *Biometrika* **83**: 329-341.

[3] Dutta, S. (2012). Local smoothing using the bootstrap. Communications in Statistics-Simulation and Computation. Accepted for publication.

[4] Efron, B. (1979). Bootstrap methods: Another look at the jackknife, *The Annals of Statistics* **7**:1-26.

[5] Efron, B. and Tibshirani, R. J. (1986). Bootstrap methods for standard error, confidence intervals, and other measures of statistical accuracy. *Statist. Science* **1** :54-77.

[6] Hall, P. (1992). Effect of Bias Estimation on Coverage Accuracy of Bootstrap Confidence Intervals for a Probability Density. *The Annals of Statistics* **20**, 2 : 675-694.

[7] Horowitz, J. L. (2001). The bootstrap. In Handbook of Econometrics, ed. J. J. Heckman and E. Leamer, vol. 5, 31593228. Amsterdam: North-Holland.

[8] Fiorio, C. V. (2004). Confidence intervals for kernel density estimation. *The Stata Journal* **4**, 2: 168179.

[9] Sheather, S.J.,and Jones, M.C. (1991). A Reliable Data-Based Bandwidth Selection Method for Kernel Density Estimation. J. Roy. *Statist. Soc. Ser. B*. **53** : 683-690.

[10] Shao, J. and Tu, D. (1995). The Jackknife and Bootstrap. Springer-Verlag, New-York.

Santanu Dutta
Mathematical Science Department
Tezpur University
Napaam:784028, Tezpur, Assam, INDIA,
tezpur1976@gmail.com