# Adaptive Histograms from a Randomized Queue that is Prioritized for Statistically Equivalent Blocks

### Gloria Teng Jennifer Harlow Raazesh Sainudiin

Department of Mathematics and Statistics, University of Canterbury, New Zealand

August 19, 2010

・ 同 ト ・ ヨ ト ・ ヨ

### Introduction

- Present statistical regular sub-pavings as an efficient, data-driven, multi-dimensional data-structure for non-parametric density estimation of massive data sets;
- Apply our methods to earthquakes in NZ, weather and aircraft trajectories over a busy US airport and samples simulated from challenging multi-dimensional densities, including Levy and Rosenbrock.



Figure: Shape of a Levy density with 700 modes.

Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

### Intervals and Boxes in $\mathbb{R}^{d}$

Intervals and Boxes as interval vectors:

$$\mathbf{x} = [\underline{x}_1, \overline{x}_1] \times [\underline{x}_2, \overline{x}_2] \times \ldots \times [\underline{x}_d, \overline{x}_d], \, \underline{x}_i \leq \overline{x}_i \;\; .$$



Figure: Boxes in 1D, 2D, and 3D.

< ロ > < 同 > < 三 > < 三

Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

## **Binary Tree Representation**

These boxes can also be represented by ordered binary trees. An operation of bisection on a box is equivalent to performing the operation on its corresponding node in the tree, i.e.:



Figure: Bisecting a box or its equivalent node.

Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

Regular Sub-pavings (RSPs) (Jaulin et. al., 2001)

- A sequence of bisections of boxes;
- Start from the root box;
- Along the first widest dimension.

Figure: A sequence of bisections on root box  $\mathbb{X}$  to produce a 4-leafed RSP *s*.



Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

Regular Sub-pavings (RSPs) (Jaulin et. al., 2001)

- A sequence of bisections of boxes;
- Start from the root box;
- Along the first widest dimension.

Figure: A sequence of bisections on root box X to produce a 4-leafed RSP *s*.



→ Ξ →

Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

Regular Sub-pavings (RSPs) (Jaulin et. al., 2001)

- A sequence of bisections of boxes;
- Start from the root box;
- Along the first widest dimension.

Figure: A sequence of bisections on root box X to produce a 4-leafed RSP *s*.



Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

Regular Sub-pavings (RSPs) (Jaulin et. al., 2001)

- A sequence of bisections of boxes;
- Start from the root box;
- Along the first widest dimension.

Figure: A sequence of bisections on root box  $\mathbb X$  to produce a 4-leafed RSP *s*.



Teng, Harlow and Sainudiin

Adaptive Histograms from SEB-based PQ

Adaptive Histograms Arithmetic on SRSPs Application Conclusion Intervals and Boxes **Regular Sub-pavings (RSPs)** Statistical Regular Sub-pavings (SRSPs)

### The Space of All Possible RSPs

The number of distinct RSP with i splits is equal to the Catalan number:



A B F A B F

Adaptive Histograms Arithmetic on SRSPs Application Conclusion Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

# Statistical Regular Sub-pavings (SRSPs)

- Extended from the RSP;
- Caches recursively computable statistics at each box or node as data falls through;
- These statistics include:
  - the sample count;
  - the sample mean vector;
  - the sample variance-covariance matrix;
  - and the volume of the box.

Figure: Caching the sample count in each node (or box).





Adaptive Histograms Arithmetic on SRSPs Application Conclusion Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

# Statistical Regular Sub-pavings (SRSPs)

- Extended from the RSP;
- Caches recursively computable statistics at each box or node as data falls through;
- These statistics include:
  - the sample count;
  - the sample mean vector;
  - the sample variance-covariance matrix;
  - and the volume of the box.

Figure: Caching the sample count in each node (or box).





Adaptive Histograms Arithmetic on SRSPs Application Conclusion Intervals and Boxes Regular Sub-pavings (RSPs) Statistical Regular Sub-pavings (SRSPs)

# Statistical Regular Sub-pavings (SRSPs)

- Extended from the RSP;
- Caches recursively computable statistics at each box or node as data falls through;
- These statistics include:
  - the sample count;
  - the sample mean vector;
  - the sample variance-covariance matrix;
  - and the volume of the box.

Figure: Caching the sample count in each node (or box).



S.E.B. Priority Queue

### SRSPs as Adaptive Histograms

The histogram estimate of i.i.d. random variables  $X_1, X_2, \ldots, X_n$  in  $\mathbb{R}^d$  with density f is given by:

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{I_{X_i \in \mathbf{X}(x)}}{vol(\mathbf{x})}$$

 $\mathbf{x}(x)$ : the leaf box  $\mathbf{x}$  that contains x

 $vol(\mathbf{x})$ : volume of box  $\mathbf{x}$ 

Figure: A SRSP as a histogram estimate.



S.E.B. Priority Queue

### SRSPs as Adaptive Histograms

The histogram estimate of i.i.d. random variables  $X_1, X_2, \ldots, X_n$  in  $\mathbb{R}^d$  with density f is given by:

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{I_{X_i \in \mathbf{X}(x)}}{vol(\mathbf{x})}$$

 $\mathbf{x}(x)$ : the leaf box  $\mathbf{x}$  that contains x

 $vol(\mathbf{x})$ : volume of box  $\mathbf{x}$ 

Figure: A SRSP as a histogram estimate.





Teng, Harlow and Sainudiin Adaptive Histograms from SEB-based PQ

S.E.B. Priority Queue

## A Prioritized Queue based Algorithm

#### Algorithm SplitMostCounts

As data arrives, order the leaf boxes of the SRSP so that the leaf box with **the most number of points** will be chosen for the next bisection.



S.E.B. Priority Queue

## A Prioritized Queue based Algorithm

#### Algorithm SplitMostCounts

As data arrives, order the leaf boxes of the SRSP so that the leaf box with **the most number of points** will be chosen for the next bisection.



S.E.B. Priority Queue

## A Prioritized Queue based Algorithm

#### Algorithm SplitMostCounts

As data arrives, order the leaf boxes of the SRSP so that the leaf box with **the most number of points** will be chosen for the next bisection.

Two or more boxes with the most number of points?



S.E.B. Priority Queue

## A Prioritized Queue based Algorithm

#### Algorithm SplitMostCounts

As data arrives, order the leaf boxes of the SRSP so that the leaf box with **the most number of points** will be chosen for the next bisection.

Break ties by picking these boxes at random for the next bisection.



S.E.B. Priority Queue

## A Prioritized Queue based Algorithm

#### Algorithm SplitMostCounts

As data arrives, order the leaf boxes of the SRSP so that the leaf box with **the most number of points** will be chosen for the next bisection.

Keep bisecting till each box has less than or equal to  $k_n$  number of points (let  $k_n = 3$  here).



S.E.B. Priority Queue

## A Prioritized Queue based Algorithm

#### Algorithm SplitMostCounts

As data arrives, order the leaf boxes of the SRSP so that the leaf box with **the most number of points** will be chosen for the next bisection.



Some Examples

Figure: Histogram density estimates their corresponding sub-pavings for the bivariate Gaussian, Levy and Rosenbrock densities.

S.E.B. Priority Queue



Teng, Harlow and Sainudiin

Adaptive Histograms from SEB-based PQ

### Choice of $k_n$



Figure: Two histogram density estimates for the standard bivariate gaussian density with different choices of  $k_n$ . The histogram is under-smoothed when  $k_n$  is relatively smaller than n and over-smoothed when  $k_n$  is relatively larger.

Teng, Harlow and Sainudiin

## Adding and Averaging SRSPs

Perform a non-minimal union (or add sub-pavings) and adjust counts:



< 同 > < 三 > < 三 >

## Adding and Averaging SRSPs

Perform a non-minimal union (or add sub-pavings) and adjust counts:



э

### Adding and Averaging SRSPs

### Adding *m* histogram density estimates

$$\sum_{i=1}^{m} \hat{f}^{(i)} = \hat{f}^{(1)} + \hat{f}^{(2)} + \hat{f}^{(3)} + \dots + \hat{f}^{(m)}$$
$$= \left( \left( \left( \hat{f}^{(1)} + \hat{f}^{(2)} \right) + \hat{f}^{(3)} \right) + \dots + \hat{f}^{(m)} \right) .$$

- 4 同 2 4 日 2 4 日 2

э

## Adding and Averaging SRSPs

### Adding *m* histogram density estimates

$$\sum_{i=1}^{m} \hat{f}^{(i)} = \hat{f}^{(1)} + \hat{f}^{(2)} + \hat{f}^{(3)} + \dots + \hat{f}^{(m)}$$
$$= \left( \left( \left( \hat{f}^{(1)} + \hat{f}^{(2)} \right) + \hat{f}^{(3)} \right) + \dots + \hat{f}^{(m)} \right) .$$

Averaging *m* histogram density estimate

$$\overline{\widehat{f}} = \frac{1}{m} \sum_{i=1}^{m} \widehat{f}^{(i)}$$

### An Example



Figure: Histogram density estimates of the bivariate Levy using different values of  $k_n$ .

### An Example



Figure: The averaged histogram density estimate.

Image: A math and A

### An Example of Application

#### Example

Air Traffic Data (Link to SAGE server): interested in applying SRSPs to the analysis of thunderstorm effects on aggregated aircraft trajectories.

### Conclusions

- We proposed an efficient, data-driven, multi-dimensional data-structure, SRSPs, for non-parametric density estimation of massive data sets;
- The SRSP can be represented by a binary tree and can either grow (through bisection of nodes) or be pruned (through merging nodes) adaptively;
- Arithmetic operations can be efficiently extended to these data structures, i.e. averaging histograms.

🗇 🕨 🔺 🖻 🕨 🔺 🖻

### References

Jaulin, L., Kieffer, M., Didrit, O. & Walter, E. (2001). *Applied interval analysis*. London: Springer-Verlag.

Lugosi, G. and Nobel, A. (1996). Consistency of data-driven histogram methods for density estimation and classification. *The Annals of Statistics* **24** 687–706.

Sainudiin, R. and York, T. L. (2005). *An Auto-validating Rejection Sampler*. BSCB Dept. Technical Report BU-1661-M, Cornell University, Ithaca, New York.

Tucker, W. (2004). *Auto-validating numerical methods*. Lecture Notes, Uppsala University, Sweden.

# Thank you!

- 4 同 2 4 日 2 4 日 2