

A Fast Approximation Scheme for Probabilistic Wavelet Synopses

Antonios Deligiannakis
University of Maryland
adel@cs.umd.edu

Minos Garofalakis
Bell Laboratories
minos@research.bell-labs.com

Nick Roussopoulos
University of Maryland
nick@cs.umd.edu

Abstract

Several studies have demonstrated the effectiveness of Haar wavelets in reducing large amounts of data down to compact wavelet synopses that can be used to obtain fast, accurate approximate query answers. While Haar wavelets were originally designed for minimizing the overall root-mean-squared (i.e., L_2 -norm) error in the data approximation, the recently-proposed idea of probabilistic wavelet synopses also enables their use in minimizing other error metrics, such as the relative error in individual data-value reconstruction, which is arguably the most important for approximate query processing. Known construction algorithms for probabilistic wavelet synopses employ probabilistic schemes for coefficient thresholding that are based on optimal Dynamic-Programming (DP) formulations over the error-tree structure for Haar coefficients. Unfortunately, these (exact) schemes can scale quite poorly for large data-domain and synopsis sizes. To address this shortcoming, in this paper, we introduce a novel, fast approximation scheme for building probabilistic wavelet synopses over large data sets. Our algorithm’s running time is near-linear in the size of the data-domain (even for very large synopsis sizes) and proportional to $1/\epsilon$, where ϵ is the desired approximation guarantee. The key technical idea in our approximation scheme is to make exact DP formulations for probabilistic thresholding much “sparser”, while ensuring a maximum relative degradation of ϵ on the quality of the approximate synopsis, i.e., the desired approximation error metric. Extensive experimental results over synthetic and real-life data clearly demonstrate the benefits of our proposed techniques.

1. Introduction

Approximate query processing over compact, pre-computed data synopses has attracted a lot of interest recently as a viable solution for dealing with complex queries over massive amounts of data in interactive decision-support and data-exploration environments. For several of these application scenarios, exact answers are not required, and users may in fact prefer fast, approximate answers to their queries. Examples include the initial, exploratory drill-down queries in ad-hoc data mining systems, where the goal is to quickly identify the “interesting” regions of the underlying database; or, aggregation

queries in decision-support systems where the full precision of the exact answer is not needed and the first few digits of precision suffice (e.g., the leading digits of a total in the millions or the nearest percentile of a percentage) [1, 2, 6, 12].

Background and Earlier Results. *Haar wavelets* are a mathematical tool for the hierarchical decomposition of functions with several successful applications in signal and image processing [13, 18]. A number of recent studies has also demonstrated the effectiveness of the Haar wavelet decomposition as a data-reduction tool for database problems, including selectivity estimation [14] and approximate query processing over massive relational tables [2, 19] and data streams [9, 15]. Briefly, the key idea is to apply the decomposition process over an input data set along with a thresholding procedure in order to obtain a compact data synopsis comprising of a selected small set of *Haar wavelet coefficients*. The results of the recent research studies of Matias, Vitter and Wang [14, 19], Chakrabarti et al. [2, 3], and others [5, 17] have demonstrated that fast and accurate approximate query processing engines can be designed to operate solely over such compact *wavelet synopses*.

Until very recently, a major criticism of wavelet-based approximate query processing techniques has been the fact that unlike, e.g., random samples, conventional wavelet synopses (such as those used in all the above-cited studies) cannot provide useful guarantees on the quality of approximate answers. The problem here is that coefficients for such conventional synopses are typically chosen in a greedy fashion in order to optimize the overall root-mean-squared (i.e., L_2 -norm) error in the data approximation. However, as pointed out by Garofalakis and Gibbons [7], conventional, L_2 -optimized wavelet synopses can result in approximate answers of widely-varying quality (even within the same data set) and approximation errors that are heavily biased towards certain regions of the underlying data domain. Their proposed solution, termed *probabilistic wavelet synopses* [7], employs the idea of *randomized coefficient rounding* in conjunction with *Dynamic-Programming-based* thresholding schemes specifically tuned for optimiz-

ing the *maximum relative error* in the approximate reconstruction of individual data values. By optimizing for relative error (with a sanity bound), which is arguably the most important metric for approximate query answers, probabilistic wavelet synopses offer drastic reductions in the approximation error over conventional deterministic techniques and, furthermore, enable unbiased data reconstruction with meaningful, non-trivial *error guarantees* for reconstructed values [7].¹

Our Contributions. The Dynamic-Programming (DP) algorithms of [7] for constructing probabilistic wavelet synopses are based on an *optimal, continuous DP formulation* over the error-tree structure for Haar coefficients, in conjunction with the idea of *quantizing* the possible choices for synopsis-space allocation using an integer parameter $\alpha > 1$ (in other words, fractional space is allotted to coefficients in multiples of $1/\alpha$). Unfortunately, the problem with these exact (modulo the quantization) DP techniques is that they can scale poorly for large data-domain and synopsis sizes – with a domain size of N and synopsis storage of B , the worst-case running time of the optimized algorithm presented in [7] (which uses binary-search to optimize the DP search) is $O(N\alpha^2 B \log(\alpha B))$, which becomes $O(N^2\alpha^2 \log(N\alpha))$ for large synopsis sizes $B = \Theta(N)$. (Given that today’s personal computers and workstations typically come equipped with Gigabytes of main memory, it is quite realistic to expect large synopsis sizes when dealing with massive data sets.) Our own experience with the DP schemes in [7] has demonstrated that the times required for building a probabilistic wavelet synopsis can increase very rapidly for large domain sizes N and synopsis sizes B ; this certainly raises some concerns with respect to the applicability of probabilistic wavelet techniques on massive, real-life data sets. Note that large domain sizes in the range of 10^5 – 10^7 are not at all uncommon, e.g., for massive time-series data sets where one or more readings/measurements are continuously recorded on every time-tick.

To address these concerns, we propose a novel, fast *approximation scheme* for building probabilistic wavelet synopses over large data sets. Given a quantization parameter α and a desired approximation factor ϵ , our algorithm can be used to build a probabilistic synopsis of *any size* $B \leq N$ in worst-case time of $O(N\alpha \log \alpha \min\{\log N \log R/\epsilon, B\alpha\})$ (where R is roughly proportional to the maximum absolute Haar-coefficient value in the decomposition), while guaranteeing that the quality of the final solution is within a factor of $(1 + \epsilon)$ of that obtained by the (exact) techniques of Garofalakis and Gibbons [7] for the same problem in-

stance. (Note that the running time of our algorithm actually represents an improvement over the techniques in [7] even when computing the exact, optimal solution.) In a nutshell, the key technical idea in our proposed approximation scheme is to make the original DP formulations in [7] *much “sparser”*, while ensuring a maximum relative degradation of $(1 + \epsilon)$ on the quality of the approximate solution, i.e., the desired maximum error metric. This is accomplished by restricting the DP search to a carefully-chosen, logarithmically-small subset of “*breakpoints*” that cover the entire range of possible space allotments within the required error guarantee. Our results clearly validate our approach, demonstrating that our algorithm (1) exhibits significantly smaller running times, often by *more than one or even two orders of magnitude*, than the exact DP solution; and (2) typically produces significantly tighter approximations than the specified $(1 + \epsilon)$ (worst-case) guarantee.

Roadmap. The remainder of this paper is organized as follows. Section 2 gives background material on wavelets, as well as conventional and probabilistic wavelet synopses. In Section 3, we discuss our approximation scheme for constructing probabilistic wavelet synopses in detail. Section 4 describes the results of our empirical study and, finally, Section 5 gives some concluding remarks.

2. Preliminaries

The Haar Wavelet Transform. Wavelets are a useful mathematical tool for hierarchically decomposing functions in ways that are both efficient and theoretically sound. Broadly speaking, the wavelet decomposition of a function consists of a coarse overall approximation along with detail coefficients that influence the function at various scales [18]. Suppose that we are given the one-dimensional data vector A containing the $N = 8$ data values $A = [2, 2, 0, 2, 3, 5, 4, 4]$. The Haar wavelet transform of A can be computed as follows. We first average the values together pairwise to get a new “lower-resolution” representation of the data with the following average values $[2, 1, 4, 4]$. In other words, the average of the first two values (that is, 2 and 2) is 2, that of the next two values (that is, 0 and 2) is 1, and so on. Obviously, some information has been lost in this averaging process. To be able to restore the original values of the data array, we need to store some *detail coefficients*, that capture the missing information. In Haar wavelets, these detail coefficients are simply the differences of the (second of the) averaged values from the computed pairwise average. Thus, in our simple example, for the first pair of averaged values, the detail coefficient is 0 since $2 - 2 = 0$, for the second we need to store -1 since $1 - 2 = -1$. Note that no information has been lost in this process – it is fairly simple to reconstruct the eight values of the original data array from the lower-resolution array containing the four av-

¹ In more recent work, Garofalakis and Kumar [8] have proposed optimal *deterministic* wavelet-thresholding schemes for relative error metrics; still, their optimal algorithms are significantly more expensive computationally than the probabilistic schemes in [7], and do not directly extend to multi-dimensional data.

verages and the four detail coefficients. Recursively applying the above pairwise averaging and differencing process on the lower-resolution array containing the averages, we get the following full decomposition:

Resolution	Averages	Detail Coefficients
3	[2, 2, 0, 2, 3, 5, 4, 4]	—
2	[2, 1, 4, 4]	[0, -1, -1, 0]
1	[3/2, 4]	[1/2, 0]
0	[11/4]	[-5/4]

The *wavelet transform* (also known as the *wavelet decomposition*) of A is the single coefficient representing the overall average of the data values followed by the detail coefficients in the order of increasing resolution. Thus, the one-dimensional Haar wavelet transform of A is given by $W_A = [11/4, -5/4, 1/2, 0, 0, -1, -1, 0]$. Each entry in W_A is called a *wavelet coefficient*. The main advantage of using W_A instead of the original data vector A is that for vectors containing similar values most of the detail coefficients tend to have very small values. Thus, eliminating such small coefficients from the wavelet transform (i.e., treating them as zeros) introduces only small errors when reconstructing the original data, resulting in a very effective form of lossy data compression [18]. Furthermore, the Haar wavelet decomposition can also be extended to *multi-dimensional* data arrays through natural generalizations of the one-dimensional decomposition process described above. Multi-dimensional Haar wavelets have been used in a wide variety of applications, including approximate query answering over complex decision-support data sets [2, 19].

Error Tree and Conventional Wavelet Synopses. A helpful tool for exploring the properties of the Haar wavelet decomposition is the *error tree* structure [14]. The error tree is a hierarchical structure built based on the wavelet transform process. Figure 1 depicts the error tree for our example data vector A . Each internal node c_i ($i = 0, \dots, 7$) is associated with a wavelet coefficient value, and each leaf d_i ($i = 0, \dots, 7$) is associated with a value in the original data array; in both cases, the index i denotes the positions in the data array or error tree. For example, c_0 corresponds to the overall average of A . The resolution levels l for the coefficients (corresponding to levels in the tree) are also depicted. We use the terms “node” and “coefficient” interchangeably in what follows.

Given a node u in an error tree T , let $\text{path}(u)$ denote the set of all proper ancestors of u in T (i.e., the nodes on the path from u to the root of T , including the root but not u) with non-zero coefficients. A key property of the Haar wavelet decomposition is that the reconstruction of any data value d_i depends only on the values of coefficients on $\text{path}(d_i)$; more specifically, we have $d_i = \sum_{c_j \in \text{path}(d_i)} \delta_{ij} \cdot c_j$, where $\delta_{ij} = +1$ if d_i is in the left child subtree of c_j or $j = 0$, and $\delta_{ij} = -1$ otherwise. For example, in Figure 1, $d_4 = c_0 - c_1 + c_6 = \frac{11}{4} - (-\frac{5}{4}) + (-1) = 3$.

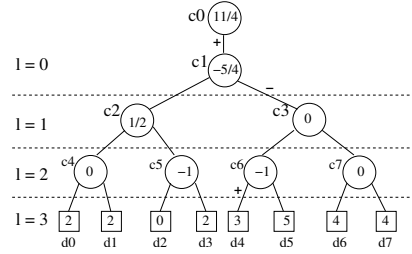


Figure 1. Error tree for example array A ($N = 8$).

Given a limited amount of storage for building a *wavelet synopsis* of the input data array A , a thresholding procedure retains a certain number $B \ll N$ of the coefficients as a highly-compressed approximate representation of the original data (the remaining coefficients are implicitly set to 0). Conventional coefficient thresholding is a deterministic process that seeks to minimize the overall root-mean-squared error (L_2 error norm) of the data approximation [18] by retaining the B largest wavelet coefficients in *absolute normalized value* [18]. L_2 coefficient thresholding has also been the method of choice for the bulk of existing work on Haar-wavelets applications in the data-reduction and approximate query processing domains [2, 14, 15, 19].

Probabilistic Wavelet Synopses. Unfortunately, wavelet synopses optimized for overall L_2 error using the above-described process may not always be the best choice for approximate query processing systems. As observed in a recent study by Garofalakis and Gibbons [7], such conventional wavelet synopses suffer from several important problems, including the introduction of severe bias in the data reconstruction and wide variance in the quality of the data approximation, as well as the lack of non-trivial guarantees for individual approximate answers. To address these shortcomings, their work introduces *probabilistic wavelet synopses*, a novel approach for constructing data summaries from wavelet-transform arrays. In a nutshell, their key idea is to apply a probabilistic thresholding process based on *randomized rounding* [16], that randomly rounds coefficients either up to a larger rounding value or down to zero, so that the value of each coefficient is correct *on expectation*. More formally, each non-zero wavelet coefficient c_i is associated with a *rounding value* λ_i and a corresponding *retention probability* $y_i = \frac{c_i}{\lambda_i}$ such that $0 < y_i \leq 1$, and the value of coefficient c_i in the synopsis becomes a random variable $C_i \in \{0, \lambda_i\}$, where,

$$C_i = \begin{cases} \lambda_i & \text{with probability } y_i \\ 0 & \text{with probability } 1 - y_i. \end{cases}$$

In other words, a probabilistic wavelet synopsis essentially “rounds” each non-zero wavelet coefficient c_i *independently* to either λ_i or zero by flipping a biased coin with

$$M^*[i, \beta_i] = \begin{cases} \min_{\substack{y_i \in (0, \min\{1, \beta_i\}] \\ b_L \in [0, \beta_i - y_i]}} \left\{ \max \left\{ \frac{\text{Var}(i, y_i)}{\text{Norm}(2i)} + M^*[2i, b_L], \right. \right. \\ \left. \left. \frac{\text{Var}(i, y_i)}{\text{Norm}(2i+1)} + M^*[2i+1, \beta_i - y_i - b_L] \right\} \right\} & \text{if } i < N, c_i \neq 0, \\ & \text{and } \beta_i > 0 \\ \min_{b_L \in [0, \beta_i]} \{ \max\{ M^*[2i, b_L], M^*[2i+1, \beta_i - b_L] \} \} & \text{if } i < N \text{ and} \\ & c_i = 0 \\ 0 & \text{if } i \geq N \\ \infty & \text{otherwise} \end{cases} \quad (1)$$

success probability y_i . Note that the above rounding process is *unbiased*; that is, the expected value of each rounded coefficient is $E[C_i] = \lambda_i \cdot y_i + 0 \cdot (1 - y_i) = c_i$, i.e., the actual coefficient value, while its variance is

$$\text{Var}(i, y_i) = \text{Var}(C_i) = (\lambda_i - c_i) \cdot c_i = \frac{1 - y_i}{y_i} \cdot c_i^2 \quad (2)$$

and the expected size of the synopsis is simply $E[\text{synopsis}] = \sum_{i|c_i \neq 0} y_i = \sum_{i|c_i \neq 0} \frac{c_i}{\lambda_i}$. Thus, since each data value can be reconstructed as a simple linear combination of wavelet coefficients, and by linearity of expectation, it is easy to see that probabilistic wavelet synopses guarantee unbiased approximations of individual data values as well as range-aggregate query answers [7].

Garofalakis and Gibbons [7] propose several different algorithms for building probabilistic wavelet synopses. The key, of course, is to select the coefficient rounding values $\{\lambda_i\}$ such that some desired error metric for the data approximation is minimized while not exceeding a prescribed space limit B for the synopsis (i.e., $E[\text{synopsis}] \leq B$). Their winning strategies are based on formulating appropriate *Dynamic-Programming (DP)* recurrences over the Haar error-tree that explicitly minimize either (a) the maximum normalized standard error (MinRelVar), or (b) the maximum normalized bias (MinRelBias), for each reconstructed value in the data domain. As explained in [7], the rationale for these probabilistic error metrics is that they are directly related to the *maximum relative error* (with an appropriate *sanity bound* s)² in the approximation of individual data values based on the synopsis; that is, both the MinRelVar and MinRelBias schemes try to (probabilistically) control the quantity $\max_i \left\{ \frac{|\hat{d}_i - d_i|}{\max\{|d_i|, s\}} \right\}$, where \hat{d}_i denotes the data value reconstructed based on the wavelet synopsis. Note, of course, that \hat{d}_i is again a *random variable*, defined as the ± 1 summation of all (independent) coefficient random variables on $\text{path}(d_i)$. Bounding the maximum relative error in the approximation also allows for meaningful *error guarantees* to be provided on reconstructed data values [7].

As an example, Equation (1) depicts the DP recurrence in [7] for minimizing the maximum squared Normalized Stan-

dard Error (NSE²) in the data reconstruction, defined as

$$\max_i \text{NSE}^2(\hat{d}_i) = \max_i \frac{\text{Var}(\hat{d}_i)}{\max\{d_i^2, s^2\}},$$

where $\text{Var}(\hat{d}_i) = \sum_{c_j \in \text{path}(d_i)} \text{Var}(j, y_j)$. $M^*[i, \beta_i]$ here denotes the minimum value of the maximum squared NSE (i.e., NSE²) among all data values in the subtree of the error-tree rooted at coefficient c_i assuming a space budget of β_i , and $\text{Norm}(i) = \max\{d_{\min}(i)^2, s^2\}$, where $d_{\min}(i)$ is the minimum absolute data value under c_i 's subtree, is a normalization term for that subtree. (Indices $2i$ and $2i+1$ in the recurrence correspond to the left and right child (respectively) of c_i in the error-tree structure (Figure 1).) Intuitively, the DP recurrence in Equation (1) states that, for a given space budget β_i at c_i , the optimal fractional-storage allotments $\{y_i\}$ and the corresponding maximum NSE² are fixed by minimizing the larger of the costs for paths via c_i 's two child subtrees (including the root in all paths), where the cost for a path via a subtree is the sum of: (1) the variance penalty incurred at c_i itself, assuming a setting of y_i , divided by the normalization term for that subtree, and (2) the optimal cost for the subtree, assuming the given space budget. This minimization, of course, is over all possible values of y_i and, given a setting of y_i , over all possible allotments of the remaining $\beta_i - y_i$ space ‘‘units’’ amongst the two child subtrees of c_i . Of course, if $c_i = 0$ then no space budget needs to be allocated to node i , which results in the simpler recurrence in the second clause of Equation (1). Finally, data-value nodes (characterized by indices $i \geq N$, see Figure 1) cost no space and incur no cost, and the ‘‘otherwise’’ clause handles the case where we have a non-zero coefficient but zero budget ($c_i \neq 0$ and $\beta_i = 0$).

As demonstrated in [7], the DP recurrence in Equation (1) characterizes the optimal solution to the maximum NSE² minimization problem for the case of *continuous* fractional-storage allotments $y_i \in (0, 1]$ (modulo certain technical conditions that may require small ‘‘perturbations’’ of zero coefficients [7]). A similar DP recurrence can also be given for the maximum normalized bias metric. Their MinRelVar and MinRelBias algorithms then proceed by *quantizing the solution space*; that is, they assume the storage allotment variables y_i and b_L in Equation (1) to take values from a discrete set of choices corresponding to integer multiples of $1/q$, where $q > 1$ is an input integer parameter to the algorithms. (For instance, $y_i \in \{0, \frac{1}{q}, \frac{2}{q}, \dots, 1\}$ – larger values of q imply results

² The role of the sanity bound is to ensure that relative-error numbers are not unduly dominated by small data values [11, 19].

closer to the optimal, continuous solution.) Furthermore, both MinRelVar and MinRelBias cap the variance of a coefficient c at c^2 , thus allowing for zero-space allotments to unimportant coefficients (this also implies that non-zero allotments of size $\leq \frac{1}{2}$ are useless, as they result in larger variance (Equation (2)) while utilizing more space). The running time of their (quantized) MinRelVar and MinRelBias algorithms is $O(N\varrho^2 B \log(\varrho B))$ with an overall space requirement of $O(N\varrho B)$ (and an in-memory working-set size of $O(\varrho B \log N)$); furthermore, their techniques also naturally extend to multi-dimensional data and wavelets, with a reasonable increase in time and space complexity [7]. Experimental results over synthetic and real-life data in [7] have demonstrated the superiority of MinRelVar and MinRelBias probabilistic synopses as an approximate query answering tool over conventional wavelet synopses. In our discussion, we use $M_{\varrho}^*[i, \beta_i]$ to denote the result of the quantized (exact) algorithms of [7] (e.g., maximum NSE^2 for MinRelVar) for the error subtree rooted at coefficient c_i assuming a space budget of β_i .

3. Our Approximation Scheme

In this section, we present our efficient approximation scheme, termed ϵ -ApproxRV, for constructing probabilistic wavelet synopses over large data sets. Our ϵ -ApproxRV is a guaranteed $(1 + \epsilon)$ approximation algorithm for the MinRelVar scheme of Garofalakis and Gibbons [7]; that is, it focuses on minimizing the maximum NSE^2 in the data reconstruction. Our techniques can easily be extended to handle other error metrics, such as the maximum normalized bias employed by MinRelBias [7]. Our presentation here focuses primarily on the case of one-dimensional Haar wavelets – the details of the extension to multiple dimensions can be found in the full paper [4].

3.1. The One-Dimensional ϵ -ApproxRV Algorithm

Consider the error-tree structure for a one-dimensional Haar wavelet decomposition, and let R denote the *maximum absolute normalized value* of any coefficient in the tree, defined as

$$R = \max_i \frac{|c_i|}{\max\{d_{\min}(i), s\}},$$

where, as previously, $d_{\min}(i)$ denotes the minimum absolute data value in the subtree of node i . (Typically, e.g., for frequency-count vectors, the denominator in the above expression is > 1 , which implies that R is in the order of the maximum absolute coefficient value.) Our ϵ -ApproxRV algorithm runs in $O(N\varrho \log \varrho \min\{\frac{\log N \log R}{\epsilon}, B\varrho\})$ time and computes an approximate solution for any synopsis space budget $B \leq N$. Note that, for large synopsis sizes ($B = \Theta(N)$), the corresponding time complexity of the exact MinRelVar algorithm is significantly higher: $O(N^2 \varrho^2 \log(N\varrho))$ [7]. Of course, our ϵ -ApproxRV algo-

rithm is actually faster than the MinRelVar algorithm even for very small synopsis sizes ($B = o(\log N)$) and when the optimal solution is sought. Again, the key idea in our ϵ -ApproxRV algorithm is to speed up the DP search by making it much “sparser”³ – in a nutshell, our approximate “sparse” DP algorithm will only search over a few possible space allotments for each error subtree, which are carefully chosen to guarantee a maximum deviation of $(1 + \epsilon)$ from the optimal solution. Our ϵ -ApproxRV algorithm proceeds in a bottom-up fashion over the input error tree – to simplify the exposition in this section, we assume that levels in the error tree are numbered bottom-up, with leaf-node coefficients at level 0 and the root (overall average) at level $\log N - 1$.

The Sparse DP Approximation Scheme. Fix a quantization parameter ϱ , and let $M_{\varrho}[v, b]$ denote the approximate maximum squared NSE (NSE^2) computed by ϵ -ApproxRV for any data value in the error subtree rooted at node v . As earlier, $M_{\varrho}^*[v, b]$ is the corresponding optimal NSE^2 value computed by MinRelVar. Note that, for any node v , the $M_{\varrho}^*[v, b]$ values are clearly *monotonically decreasing* in b ; that is, $M_{\varrho}^*[v, x] \leq M_{\varrho}^*[v, y]$ for $x > y$ [7].

For the base case, consider a leaf-node coefficient v (at level 0) – clearly, in this case

$$M_{\varrho}^*[v, b] = \frac{\text{Var}(c_v, \min\{1, b\})}{\min\{\text{Norm}(2v), \text{Norm}(2v + 1)\}}$$

i.e., the maximum normalized variance of the corresponding random variable with a success probability of b ($b \in \{0, \frac{1}{\varrho}, \frac{2}{\varrho}, \dots, 1\}$ – any $b \geq 1$ obviously results in zero normalized variance). It is easy to see that all possible values for $M_{\varrho}^*[v, b]$, for any value of b , can be computed in time $O(\varrho)$, where ϱ is the designated quantization parameter. Out of these $O(\varrho)$ variance values and possible allotments to c_v , our ϵ -ApproxRV algorithm picks a subset of allotments $b_1 > \dots > b_h$ such that: (1) for each allotment $x \in [b_i, b_{i-1}]$ we have $M_{\varrho}^*[v, b_i] \leq (1 + \epsilon)M_{\varrho}^*[v, x]$; and, (2) b_1 through b_h cover the entire possible range of space allotments to c_v , i.e., $b_1 = 1$ and $b_h = 0$. This can obviously be done in $O(\varrho)$ time by simply going over all M_{ϱ}^* values and selecting b_{i+1} as the first allotment $\leq b_i$ such that $M_{\varrho}^*[v, b_{i+1}] > (1 + \epsilon)M_{\varrho}^*[v, b_i]$. Since the maximum normalized variance for a coefficient value c_v is at most (see Section 2) $\frac{c_v^2}{\min\{\text{Norm}(2v), \text{Norm}(2v+1)\}} \leq R^2$, is easy to see that the number h of allotment “breakpoints” selected in this fashion is at most $O(\log_{1+\epsilon} R^2) = O(\frac{\log R}{\log(1+\epsilon)}) \approx O(\frac{\log R}{\epsilon})$ (for small values of $\epsilon < 1$). The approximate error values

3 Guha et al. [10] also discuss sparse DP algorithms in an entirely different context, namely building approximate V -optimal histograms over linear, time-series data; in contrast, our solution focuses on Haar wavelets and works over the hierarchical error-tree structure of the wavelet decomposition.

determined by our ϵ -ApproxRV algorithm for coefficient c_v are defined only by these h breakpoints b_1, \dots, b_h – specifically, $M_{\mathcal{Q}}[v, b_i] = M_{\mathcal{Q}}^*[v, b_i] = \frac{\text{Var}(c_v, b_i)}{\min\{\text{Norm}(2v), \text{Norm}(2v+1)\}}$ for $i = 1, \dots, h$, and for any other possible allotment $x \in [b_i, b_{i-1})$, we define $M_{\mathcal{Q}}[v, x] = M_{\mathcal{Q}}[v, b_i]$. Thus, it is easy to see that, by construction, the approximation error of our ϵ -ApproxRV algorithm is bounded by a factor of $(1 + \epsilon)$ at leaf coefficients (at level 0); in other words, all dropped allotments are “covered” by a logarithmic number of breakpoints to within a $(1 + \epsilon)$ factor.

Now, proceeding inductively, consider an internal error-tree node v at level j , with children u and w (at level $j - 1$), and assume that the subtree rooted at u (w) has determined a collection of l_u (resp., l_w) error-function breakpoints $a_1 > \dots > a_{l_u}$ (resp., $b_1 > \dots > b_{l_w}$), and corresponding approximate NSE² values $M_{\mathcal{Q}}[\cdot]$, that cover the range of allotments to each subtree and such that, for each $x \in [a_i, a_{i-1})$ ($i = 2, \dots, l_u$), we have $M_{\mathcal{Q}}[u, a_i] \leq (1 + \epsilon)^j M_{\mathcal{Q}}^*[u, x]$ (and similarly for w). Our ϵ -ApproxRV algorithm computes the allotment breakpoints and approximate error values $M_{\mathcal{Q}}[\cdot]$ at the parent node v by iterating over all possible space allotments to node v and the breakpoints determined by the u and w subtrees (rather than all possible allotments to child subtrees), and retaining the minimum $M_{\mathcal{Q}}$ values for each total allotment. The following lemma shows that, for each fixed space allotment to the coefficient at node v , it actually suffices to look at only $l_u + l_w$ combinations (a_i, b_k) for the subtree allotments rather than all possible $l_u \cdot l_w$ combinations.

Lemma 1: *When minimizing the maximum (approximate) NSE² error at node v , for any fixed space allotment to node v , it suffices to consider only $l_u + l_w$ combinations of allotments (a_i, b_k) to the child subtrees rooted at u, w .* ■

Proof: Assume a fixed space allotment to the coefficient at node v , and let *leftVar* (*rightVar*) denote the variance of node v (for the given allotment) divided by the normalization factor of its left (resp., right) subtree. Let L_u denote the sorted list of approximate NSE² values $M'_{\mathcal{Q}}[u, a_i] = M_{\mathcal{Q}}[u, a_i] + \text{leftVar}$, i.e., $M_{\mathcal{Q}}[u, a_1] + \text{leftVar} < \dots < M_{\mathcal{Q}}[u, a_{l_u}] + \text{leftVar}$, with L_w defined similarly using the *rightVar* quantity and the $M_{\mathcal{Q}}[w, b_k]$ entries. Let $L = \text{merge}(L_u, L_w)$, i.e., $M'_{\mathcal{Q}}[y_1] \leq \dots \leq M'_{\mathcal{Q}}[y_{l_u+l_w}]$, where $y_i \in \{(u, a_k) : k = 1, \dots, l_u\} \cup \{(w, b_k) : k = 1, \dots, l_w\}$. Now assume that a_i space is allocated to the u -subtree of v . Then, it is easy to see that, when considering the allotment to the w -subtree, out of all the b -values that lie to the left of a_i in L we really only need to consider the rightmost b -value, say b_k – the reason of course is that lower values of $M'_{\mathcal{Q}}[w, b]$ (i.e., allotments $b > b_k$) result in configurations that use more total space without improving the error at v (since that is dominated by the u -subtree). These configurations are clearly useless in our error-minimization

procedure. For the b -values to the right of a_i in L , a similar argument again applies: when a value b_k is assumed, only the closest a -value to its left in L needs to be considered. ■

Thus, our approximate error-minimization procedure at v only needs to consider, for each fixed space allotment $s = 0, 1/\mathfrak{q}, \dots, 1$ to node v , $l_u + l_w$ breakpoint combinations (a_i, b_k) for the u and w subtrees, which can be determined easily in $O(l_u + l_w)$ time based on the proof of Lemma 1. Let $\mathcal{S}(s)$ denote the list of the obtained $l_u + l_w$ (a_i, b_k) combinations for each space allotment s to node v . The sorted list of approximate error values at node v can be computed in $O(\mathfrak{q}(l_u + l_w) \log \mathfrak{q})$ time by merging these lists using a heap structure or, alternatively, pairwise merging them in $\log \mathfrak{q}$ steps. Thus, an initial list of $O(\mathfrak{q}(l_u + l_w))$ breakpoints for the v subtree is determined based on the “useful” space-allocation configurations found through the above lemma – clearly, configurations that give the same (or, larger) NSE² values for the same (or, larger) amount of total space are useless and should be discarded. In other words, we define the initial set of space-allotment breakpoints for the v subtree as $\mathcal{C} = \{c = a_i + b_k + s : M_{\mathcal{Q}}[v, a_i + b_k + s] \leq M_{\mathcal{Q}}[v, a + b + t] \text{ for all } a + b + t \leq a_i + b_k + s\}$. (Useless configurations and configurations with space larger than B can easily be discarded in the merging pass for the $O(\mathfrak{q})$ sub-lists described above.) It is easy to verify that, based on our inductive assumption, this set of breakpoints \mathcal{C} covers the entire range of possible allotments for the v subtree; furthermore, the following lemma shows that it also preserves the approximation properties guaranteed by the individual subtrees.

Lemma 2: *Let $s_1 > \dots > s_h$ denote the sorted list of space-allotment breakpoints \mathcal{C} for the v subtree, computed as described above, and let $x \in [s_i, s_{i-1})$ for any i . Then, $M_{\mathcal{Q}}[v, s_i] \leq (1 + \epsilon)^j M_{\mathcal{Q}}^*[v, x]$, where j denotes the level of node v .* ■

Proof: Assume that the optimal error value $M_{\mathcal{Q}}^*[v, x]$ is obtained through the allotment configuration (y_u, y_w, s) , that is:

$$M_{\mathcal{Q}}^*[v, x] = \max \left\{ \begin{array}{l} \frac{\text{Var}(c_v, s)}{\text{Norm}(u)} + M_{\mathcal{Q}}^*[u, y_u] \\ \frac{\text{Var}(c_v, s)}{\text{Norm}(w)} + M_{\mathcal{Q}}^*[w, y_w] \end{array} \right.$$

where, of course, $x \geq y_u + y_w + s$. Since the breakpoints for the u and w subtrees cover all possible allotments, let $y_u \in [a_i, a_{i-1})$ and $y_w \in [b_k, b_{k-1})$. By our inductive hypothesis, it is easy to see that the configuration (a_i, b_k, s) (which is obviously examined by the ϵ -ApproxRV algorithm) will give $M_{\mathcal{Q}}[v, a_i + b_k + s] \leq (1 + \epsilon)^j M_{\mathcal{Q}}^*[v, x]$ and, clearly, $a_i + b_k + s \leq s_i \leq x$. Since $M_{\mathcal{Q}}[v, s_i] \leq M_{\mathcal{Q}}[v, a_i + b_k + s]$, the result follows. ■

A potential problem with our approximation scheme, as described so far, is that the list of space-allotment breakpoints \mathcal{C} would appear to grow exponentially as the DP moves up the error-tree levels. (So, starting with r breakpoints at the leaf nodes, we get $O(\alpha^j 2^j r)$ breakpoints at level j of the tree.) However, not all s_i 's in \mathcal{C} are necessary – we can actually “trim” \mathcal{C} to a small number of breakpoints, while incurring an additional $(1 + \epsilon)$ worst-case factor degradation on our approximation error. We perform this trimming process at every node of the error tree (except for the final, root node). More specifically, assume a chain of computed breakpoints $s_{i-k} > s_{i-k+1} > \dots > s_i$ such that, for each $l = i - k, \dots, i - 1$ we have $M_{\alpha}[v, s_i] \leq (1 + \epsilon)M_{\alpha}[v, s_l]$. Then, clearly, $M_{\alpha}[v, s_i]$ can “cover” all the points that are covered by s_{i-k}, \dots, s_{i-1} at an additional $(1 + \epsilon)$ degradation, since, for any $l = i - k, \dots, i - 1$:

$$\begin{aligned} M_{\alpha}[v, s_i] &\leq (1 + \epsilon)M_{\alpha}[v, s_l] \\ &\leq (1 + \epsilon)^{j+1}M_{\alpha}^*[v, x] \quad \forall x \in [s_l, s_{l-1}]. \end{aligned}$$

Thus, in this situation, the allotment points s_{i-k}, \dots, s_{i-1} can be eliminated and s_i can cover their ranges to within a $(1 + \epsilon)^{j+1}$ factor. Now, note that the maximum value of the overall NSE^2 value at level j (and, thus, the range of values for the $M_{\alpha}[\cdot]$ array) is certainly upper-bounded by $(j+1)R^2$. This means that the total number of breakpoints obtained in the manner described above is at most $\log_{1+\epsilon}((j+1)R^2) \approx O(\frac{\log(j+1) + \log R}{\epsilon})$, which is an upper bound for the size of our breakpoint-list constructed at level j of the error tree. Thus, with an overall computational effort of:

$$\begin{aligned} &\sum_{j=0}^{\log N - 1} O\left(\frac{N}{2^{j+1}} \frac{\alpha \log \alpha (\log(j+1) + \log R)}{\epsilon}\right) \\ &\leq O\left(\frac{N\alpha \log \alpha \log R}{\epsilon} \sum_{j=0}^{\log N - 1} \frac{1}{2^{j+1}}\right) + O\left(\frac{N\alpha \log \alpha}{\epsilon} \sum_{j=0}^{\log N - 1} \frac{j+1}{2^{j+1}}\right) \\ &= O\left(\frac{N\alpha \log \alpha \log R}{\epsilon}\right), \end{aligned}$$

we get a (collection of) approximate solutions at the root node of the error tree that are guaranteed to cover the optimal MinRelVar solutions to within a $(1 + \epsilon)^{\log N}$ factor. Then, it is easy to see that, setting $\epsilon' = \epsilon / \log N$, we get a guaranteed $(1 + \epsilon)$ approximation in time $O(\frac{N\alpha \log \alpha \log N \log R}{\epsilon'})$. Note that, to find the approximate solution for any specific choice of the allotment space B , we simply start out at the root and re-trace the steps of the algorithm for the largest root breakpoint s_i that is $\leq B$; to do that, we just need to keep track of the (a_i, b_k, s) configuration that generated each of the breakpoints at each error tree node and proceed recursively down the tree. It is easy to verify that the overall space required by the ϵ -ApproxRV algorithm is

$$\sum_{j=0}^{\log N - 1} O\left(\frac{N}{2^{j+1}} \frac{\log(j+1) + \log R}{\epsilon'}\right) = O\left(\frac{N \log N \log R}{\epsilon}\right),$$

while the working set size (maximum amount of memory-resident data) is only $O(\frac{\alpha(\log \log N + \log R)}{\epsilon'}) = O(\frac{\alpha \log N (\log \log N + \log R)}{\epsilon})$. To see this, note that ϵ -ApproxRV works in a bottom-up fashion and, when computing the breakpoint-list for a given node v , we only need access to: (1) the $l_u + l_w$ breakpoints of its child nodes in the error tree; and, (2) the $O(\alpha(l_u + l_w))$ initial breakpoints for node v that are computed just before the trimming process. Thus, the maximum working set will occur in the top-level of the error tree (level $\log N - 1$), where $l_u + l_w = O(\frac{\log R + \log \log N}{\epsilon'})$. Finally, note that, given a space budget B , we cannot have more than αB different breakpoints at any error-tree node; in other words, the size of the breakpoint list at any level- j node is at most $O(\min\{\frac{\log(j+1) + \log R}{\epsilon}, \alpha B\})$. This easily implies that the overall space required by our ϵ -ApproxRV algorithm can never exceed the space requirements of MinRelVar (i.e., $O(N\alpha B)$). Similarly, the list of (at most αB) breakpoints at each node can be computed in time $O(\alpha \times (\alpha B) \log \alpha) = O(B\alpha^2 \log \alpha)$; thus, the overall running time complexity is also upper-bounded by $O(NB\alpha^2 \log \alpha)$, giving an improvement over MinRelVar , even for very small values of B . The following theorem summarizes the results of our analysis for the one-dimensional ϵ -ApproxRV algorithm.

Theorem 3: *The ϵ -ApproxRV algorithm correctly computes a list of breakpoints at the root node such that, for any space budget $B \leq N$ and approximation factor ϵ , the estimated maximum NSE^2 value is within a factor of $(1 + \epsilon)$ from the optimal MinRelVar solution. The overall ϵ -ApproxRV computation requires $O(N\alpha \log \alpha \min\{\frac{\log N \log R}{\epsilon}, \alpha B\})$ time and $O(N \min\{\frac{\log N \log R}{\epsilon}, \alpha B\})$ space, with a working set size of $O(\alpha \min\{\log N \frac{\log \log N + \log R}{\epsilon}, B\})$. ■*

It is important to note that the above (worst-case) running-time and space complexities of our ϵ -ApproxRV algorithm are based on a pathological case where all the produced coefficients have an absolute normalized value of R . However, in most real-life data sets the wavelet decomposition process produces few coefficients of large magnitude, while the remaining coefficient values are significantly smaller. This, in turn, implies that, for most error-tree nodes, the maximum value of the overall NSE^2 value at level j will be significantly smaller than $(j+1)R^2$. This results not only in reduced space requirements for ϵ -ApproxRV (smaller breakpoint-lists stored at each node), but also in reduced running times (smaller breakpoint-lists scanned and merged). Our experimental results in Section 4 clearly validate our claims, with ϵ -ApproxRV demonstrating consistent and very significant gains over the exact MinRelVar scheme for a wide

range of input parameters and data sets.

Optimizations and Extensions. For very large data sets, it is possible that the breakpoint-lists produced by our ϵ -ApproxRV algorithm may not fit in main memory, resulting in substantial I/O during the algorithm’s execution. In the full paper[4], we propose a simple optimization to address this concern. Briefly, the key idea is to compute the breakpoint-lists for error-tree nodes in one pass using a postorder traversal, with a working set size of $O(\min\{\frac{(\log N)^2 (\log R + \log \log N)}{\epsilon}, \alpha B \log N\})$.

The full paper[4] also discusses in detail the extension of our ϵ -ApproxRV algorithm for multi-dimensional data sets. For the case of D -dimensional data, the running time of ϵ -ApproxRV becomes $O(\min\{\frac{N_z \alpha^4 \log M_{max} (\log \alpha + D) (D + \log R + \log \log M_{max})}{\epsilon}, N_z \alpha 2^D \times (\alpha B + D 2^D)\})$, where N_z denotes the number of nodes in the error tree that contain at least one non-zero coefficient, and M_{max} is the maximum domain size among all dimensions. The corresponding space requirements are $O(\min\{\frac{4^D N_z \log M_{max} (D + \log R + \log \log M_{max})}{\epsilon}, \alpha N_z B\})$. Note, of course, that in most real-life scenarios employing wavelet-based data reduction, the number of dimensions D is typically a small constant (e.g., 2–5) [2, 6, 7].

4. Experimental Study

In this section, we present an extensive experimental study of our proposed ϵ -ApproxRV algorithm for constructing probabilistic wavelet synopses over large data sets. The objective of this study is to evaluate both the scalability and the obtained accuracy of our proposed ϵ -ApproxRV algorithm when compared to the dynamic programming algorithm MinRelVar of [7] for a large variety of real-life and synthetic data sets. For the later DP solution, we used the significantly faster version of the algorithm that was very recently proposed in [7]. The main findings of our study for the ϵ -ApproxRV algorithm include:

- **Near Optimal Results.** The ϵ -ApproxRV algorithm consistently provides near-optimal solutions. Moreover, the actual deviation of the ϵ -ApproxRV solution from the optimal one is typically significantly smaller (usually by a factor larger than 5) than the specified ϵ value.
- **Significantly Faster Solution.** Our ϵ -ApproxRV algorithm provides a fast and scalable solution for constructing probabilistic synopses over large data sets. Compared to the MinRelVar algorithm of [7], the running time of the ϵ -ApproxRV algorithm is often more than an order of magnitude (and in some times more than two orders of magnitude) smaller, while at the same time providing highly-accurate answers. In fact, the ϵ -ApproxRV algorithm is significantly faster even for cases when the optimal solution is required ($\epsilon = 0$).

4.1. Testbed and Methodology

Techniques and Parameter Settings. Our experimental study compares the ϵ -ApproxRV and MinRelVar algorithms for constructing probabilistic wavelet synopses. Both algorithms utilize the quantization parameter q , which is assigned a value of 10, as suggested in [7], in our experiments. Larger values of this quantization parameter improved the running time performance of the ϵ -ApproxRV algorithm when compared to the MinRelVar algorithm, as expected by the running time complexities of the two algorithms. Finally, the sanity bound of each data set is set to its 5%-quantile data value.

Data Sets. We experiment with several one-dimensional synthetic and real-life data set. Due to space constraints we only present here the results for the real-life data sets (the performance of the algorithms in the synthetic data sets is qualitatively similar). The `Weather` data set contains meteorological measurements obtained by a station at the university of Washington (www-k12-atmos.washington.edu/k12/grayskies). This is a one-dimensional data set for which we extracted the following 6 measured quantities: wind speed, wind peak, solar irradiance, relative humidity, air temperature and dew-point temperature, and present here the results for the wind speed (`AirSpeed`) and the air temperature (`AirTemp`), which represent a noisy and a smooth signal, correspondingly. The `Phone` data set includes the total number of long distance calls per minute originating from several states in USA. We here present the results for the states of New York (NY) and Indiana (IN), with NY having large numbers of calls per minute and IN being a state with significantly fewer calls. The presented results for each real-life data set are also indicative of the results for the other measured quantities in these data sets.

Approximation Error Metrics. To compare the accuracy of the studied algorithms we focus on the maximum relative error of the approximation, since it can provide guaranteed error-bounds for the reconstruction of any individual data value. Since the objective function that both studied algorithms try to minimize is the maximum NSE^2 of any data value, for a more direct and clear comparison we present the results for this metric and for both algorithms. The results for the maximum relative error are qualitatively similar to the presented ones.

4.2. Experimental Results

Sensitivity to ϵ . We now evaluate the accuracy and running time of the ϵ -ApproxRV algorithm in comparison to the MinRelVar algorithm, using the real-life data sets. In Figures 2 and 3 we plot the running times for the two algorithms and for the two data sets, correspondingly, as we vary the value

of ϵ from 0 to 0.3. We set the domain size for all data sets to 65536, and the synopsis space to 5% of the input size. The ϵ -ApproxRV algorithm is consistently faster than the MinRelVar algorithm in both real-life data sets, often by more than an order of magnitude, and is considerably faster even when the optimal solution is required ($\epsilon = 0$). Unlike the MinRelVar algorithm which may perform multiple lookups of each computed entry,⁴ the ϵ -ApproxRV algorithm processes all node entries in a single pass, therefore resulting in significantly faster running times. We also observe in these figures that, with the increase of ϵ , the running time of ϵ -ApproxRV decreases, as the algorithm effectively prunes a larger number of breakpoints.

The corresponding NSE² values for both algorithms are presented in Figures 4 and 5. In order for the reader to be able to observe the difference in the accuracy of the two algorithms, in each figure we plot the ratio of the maximum NSE² values produced by the ϵ -ApproxRV algorithm to the corresponding results of the MinRelVar algorithm. The ϵ -ApproxRV algorithm, as expected, always provides solutions that are within the specified error factor ϵ from the optimal solution. It is interesting to note though that in all cases the produced solution is significantly closer to the optimal one (by more than a factor of 5), than the specified ϵ value. This is not surprising, as ϵ represents a worst-case error bound.

Sensitivity to the Domain Size. We now evaluate the accuracy and running time of the ϵ -ApproxRV algorithm in comparison to the MinRelVar algorithm, using the real-life data sets, when we vary the domain size of the data sets from 128 to 65536, and plot the resulting running times for the two algorithms in Figures 6 and 7. The synopsis space is set to 5% of the input size, while the value of ϵ is set to 0.10. Again, the ϵ -ApproxRV algorithm significantly outperforms the MinRelVar algorithm, with the savings in running time increasing rapidly as the domain size increases. For large domain sizes, the ϵ -ApproxRV algorithm is up to 23.8 times faster than the MinRelVar algorithm.

In Figures 8 and 9 we plot the corresponding ratios of the maximum NSE² values obtained by the two algorithms. Again, the ϵ -ApproxRV algorithm always produced solutions that are significantly closer to the optimal solution (less than 1.7% and 1.4% difference, correspondingly, for the two data sets), than the specified error factor ϵ .

Sensitivity to the Synopsis Space. In Figures 10 and 11 we present the running times for both algorithms and for the real-life data sets, as the synopsis space is varied from 1% to 30% of the size of the input. The domain size is set to 65536, while the value of ϵ is set to 0.10. The running time

of the MinRelVar algorithm increases rapidly with the increase in the used synopsis space, while the corresponding running time of the ϵ -ApproxRV algorithm remains practically unaffected. For large synopsis spaces, the ϵ -ApproxRV algorithm is more than two orders of magnitude faster than the MinRelVar algorithm. However, the solutions obtained from the ϵ -ApproxRV algorithm are again very close to the optimal ones (less than 1.7% and 1.5% difference for the two data sets), as we can see in Figures 12 and 13.

5. Conclusions

We have proposed a novel, fast approximation scheme for constructing probabilistic wavelet synopses over large data sets. Our proposed techniques employ a much “sparser” version of previously proposed Dynamic-Programming (DP) solutions, which restricts its search to a carefully chosen, logarithmically-small subset of “breakpoints” that cover the entire range of possible space allotments, while always ensuring a maximum relative degradation of $(1 + \epsilon)$ in the quality of the obtained solution. Our experimental evaluation has demonstrated that our approximation algorithm typically provides significantly tighter solutions than the maximum $(1 + \epsilon)$ error factor, while at the same time providing running times that are up to two orders of magnitude smaller than known exact DP solutions.

References

- [1] S. Acharya, P.B. Gibbons, V. Poosala, and S. Ramaswamy. “Join Synopses for Approximate Query Answering”. In *ACM SIGMOD*, 1999.
- [2] K. Chakrabarti, M. Garofalakis, R. Rastogi, and K. Shim. “Approximate Query Processing Using Wavelets”. In *VLDB*, 2000.
- [3] K. Chakrabarti, M. Garofalakis, R. Rastogi, and K. Shim. “Approximate query processing using wavelets”. *The VLDB Journal*, 10(2-3):199–223, September 2001.
- [4] A. Deligiannakis, M. Garofalakis and N. Roussopoulos. “A Fast Approximation Scheme for Probabilistic Wavelet Synopses”. Technical Report UMCP-CSD-CS TR #4643, 2005.
- [5] A. Deligiannakis and N. Roussopoulos. “Extended Wavelets for Multiple Measures”. In *ACM SIGMOD*, 2003.
- [6] M. Garofalakis and P.B. Gibbons. “Approximate Query Processing: Taming the Terabytes”. Tutorial in *VLDB*, 2001.
- [7] M. Garofalakis and P.B. Gibbons. “Probabilistic Wavelet Synopses”. *ACM Transactions on Database Systems*, 29(1), March 2004.
- [8] M. Garofalakis and A. Kumar. “Deterministic Wavelet Thresholding for Maximum-Error Metrics”. In *PODS*, 2004.
- [9] A.C. Gilbert, Y. Kotidis, S. Muthukrishnan, and M.J. Strauss. “Surfing Wavelets on Streams: One-pass Summaries for Approximate Aggregate Queries”. In *VLDB*, 2001.
- [10] S. Guha, N. Koudas and K. Shim. “Data-Streams and Histograms”. In *STOC*, 2001.
- [11] P.J. Haas and A.N. Swami. “Sequential Sampling Procedures for Query Size Estimation”. In *ACM SIGMOD*, 1992.
- [12] J.M. Hellerstein, P.J. Haas, and H.J. Wang. “Online Aggregation”. In *ACM SIGMOD*, 1997.
- [13] B. Jawerth and W. Sweldens. “An Overview of Wavelet Based Multiresolution Analyses”. *SIAM Review*, 36(3):377–412, 1994.
- [14] Y. Matias, J.S. Vitter, and M. Wang. “Wavelet-Based Histograms for Selectivity Estimation”. In *ACM SIGMOD*, 1998.

⁴ In the MinRelVar algorithm, the optimal solution of allocating space B to any node v may be probed for any space allotment $\geq B$ to v 's parent node in the error tree.

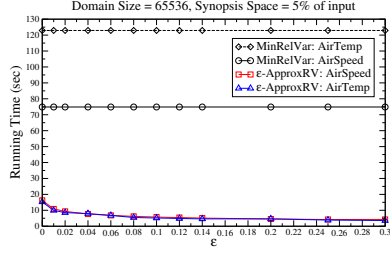


Figure 2: Running Time vs ϵ in Weather data

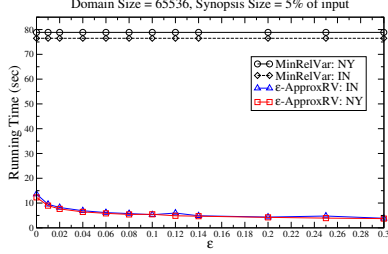


Figure 3: Running Time vs ϵ in Phone data

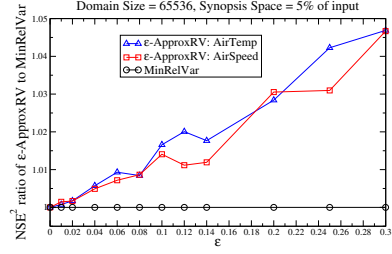


Figure 4: NSE^2 ratio vs ϵ in Weather data

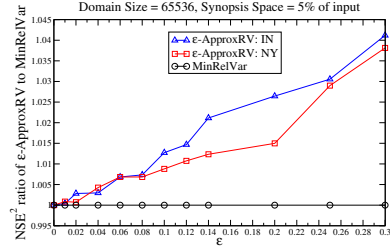


Figure 5: NSE^2 ratio vs ϵ in Phone data

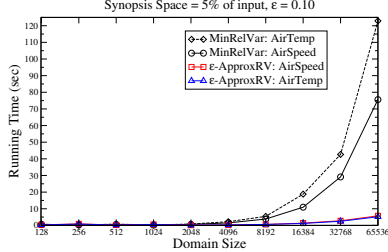


Figure 6: Running Time vs Domain Size in Weather data

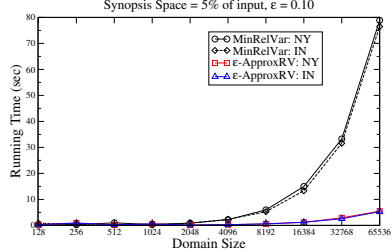


Figure 7: Running Time vs Domain Size in Phone data

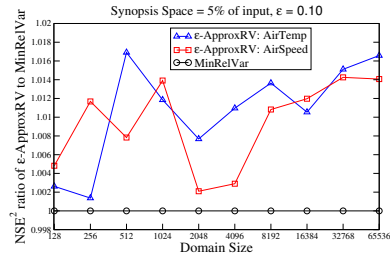


Figure 8: NSE^2 ratio vs Domain Size in Weather data

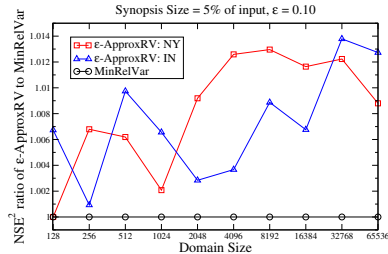


Figure 9: NSE^2 ratio vs Domain Size in Phone data

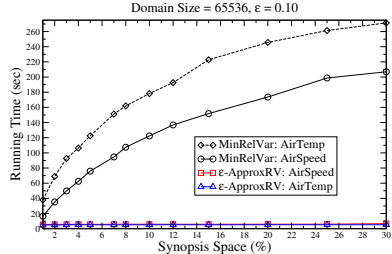


Figure 10: Running Time vs Space in Weather data

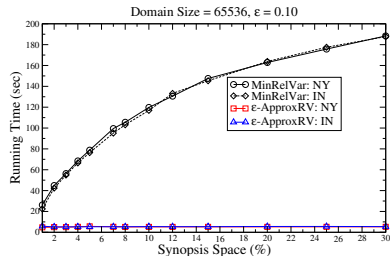


Figure 11: Running Time vs Space in Phone data

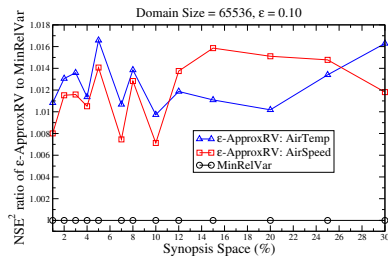


Figure 12: NSE^2 ratio vs Space in Weather data

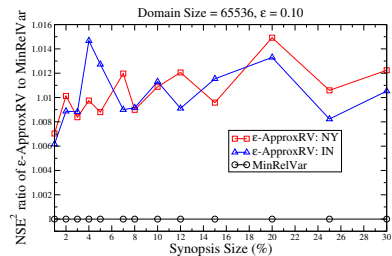


Figure 13: NSE^2 ratio vs Space in Phone data

[15] Y. Matias, J.S. Vitter, and M. Wang. "Dynamic Maintenance of Wavelet-Based Histograms". In *VLDB*, 2000.
 [16] R. Motwani and P. Raghavan. "Randomized Algorithms". Cambridge University Press, 1995.
 [17] R.R. Schmidt and C. Shahabi. "ProPolyn: A Fast Wavelet-Based Algorithm for Progressive Evaluation of Polynomial Range-Sum

Queries". In *EDBT*, 2002.
 [18] E.J. Stollnitz, T.D. DeRose, and D.H. Salesin. "Wavelets for Computer Graphics – Theory and Applications". Morgan Kaufmann Publishers, 1996.
 [19] J.S. Vitter and M. Wang. "Approximate Computation of Multidimensional Aggregates of Sparse Data Using Wavelets". In *ACM SIGMOD*, 1999.