

# XStream

## 1 Histogram Forests (HF)

**Fitting.** For each histogram  $i = 1, \dots, T$ :

1. **Initialize the random projection matrix**  $A_i$ . This may be (i) sparse with 2/3 zeros, or (ii) Gaussian random projections (which demonstrate better performance, but are computationally expensive).
2. **Project the data.**  $Y = XA_i^T$ .
3. **Sample a shift.**  $b_i \sim \text{Uniform}(0, w)$ .
4. **Bin the data.** For each projected data point  $Y_j = (y_{j1}, \dots, y_{jK})$ , its bin vector is given by  $\tilde{Y}_j = (\lfloor \frac{y_{j1} + b_i}{w} \rfloor, \dots, \lfloor \frac{y_{jK} + b_i}{w} \rfloor)$ .
5. **Store bin counts.** In an exact hashtable or count-min sketch  $H_i$ , increment the count  $H_i(\tilde{Y}_j)$  for  $j = 1, \dots, N$ .

**Scoring.** The anomaly score of a point  $X_j$  is the negative of the average of  $H_i(\tilde{Y}_j)$  across all histograms  $i$ .

### 1.1 HF Results

Methods compared in Figure 1 (on the synthetic data with noise):

- **HF:** Histogram forest with varying bin width  $w$ , number of projected dimensions  $K$  and number of histograms  $T$ .
- **IF:** iForest with 100 trees.
- **IF-P:** iForest with 100 trees and data projected to  $k$  dimensions using sparse random projections, with varying  $k$ .

Observations from Figure 1:

1. HF with  $K = 25, T = 100, w = 4.0$  is on par with iForest with 100 trees on the data projected to  $k = 50$  or  $k = 75$  dimensions, and better than iForest for  $k = 25$ .
2. HF with  $K = 50, T = 100$  has poor performance for all  $w$ . This was our configuration setting for the previous histogram experiments (with some minor variations): we considered them a failure without trying a *lower* value of  $K$ . This suggests that  $K$  as a measure of data approximation quality may not be the right interpretation.
3. HF with  $K = 50, T = 200$  improves performance over  $T = 100$  slightly, for all  $w$ .

Histogram Forest Parameters	
$K$	Dimension of projected data
$T$	Number of histograms
$w$	Bin width
Other variables	
$X$	$N \times D$ data matrix
$Y$	$N \times K$ projected data matrix
$A_i$	$K \times D$ projection matrix
$b_i$	Shift $\sim \text{Uniform}(0, w)$
$H_i$	Hashtable/count-min sketch

Table 1: Notation.

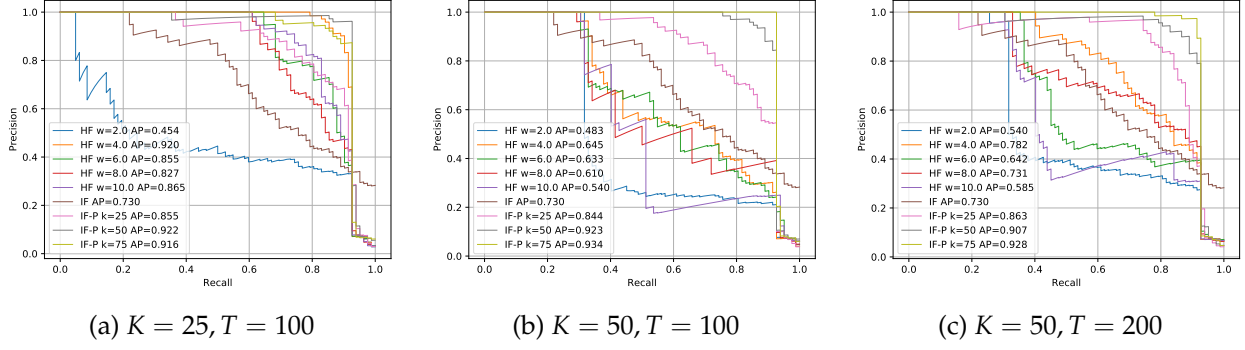


Figure 1: Precision-recall curves on synthetic data with 100 noisy dimensions.

## 1.2 Interpretation — Approximate Nearest Neighbors

The binned data point in our histogram is identical to a point “hashed” using 2-stable distributions (Datar et al. 2004, 3.1), which extends LSH to solve the approximate nearest neighbors problem in  $L_2$ . In this method, each point  $x$  is “hashed” using the function  $h(x) = \lfloor \frac{x^T a + b}{w} \rfloor$  where  $a$  is a random Gaussian vector,  $b$  is a random shift and  $w$  is the bin width. The probability that two points  $x$  and  $y$  hash to the same value is proportional to  $\|x - y\|_2$ .

Each point is hashed  $K$  times, and points that agree on all  $K$  values are considered candidate nearest neighbors. However, as  $K$  grows, it becomes less likely that even nearby points agree on all  $K$  values; hence, the process is repeated  $T$  times, and the points agreeing on all  $K$  values at in least one of the  $T$  trials are considered candidate nearest neighbors.

Each configuration of  $K, T$  and  $w$  results in a certain *probability gap* for points to be considered nearest neighbors. Specifically, this probability gap is defined by the 4-tuple  $(R, p_1, cR, p_2)$ :

- The probability of two points with distance  $\leq R$  being hashed to the same values is  $\geq p_1$ .
- The probability of two points with distance  $\geq cR$  being hashed to the same values is  $\leq p_2$ .

Hence, our anomaly score for a point can be viewed as the negative of the approximate number of points within a sphere of some radius  $R$  centered at that point. The optimal  $R$  for anomaly detection on a given dataset is data-dependent (Fig. 2). Hence, the optimal  $K, T, w$  are also data-dependent. This explains why a  $K = 25$  histogram forest may perform better than a  $K = 50$  one.

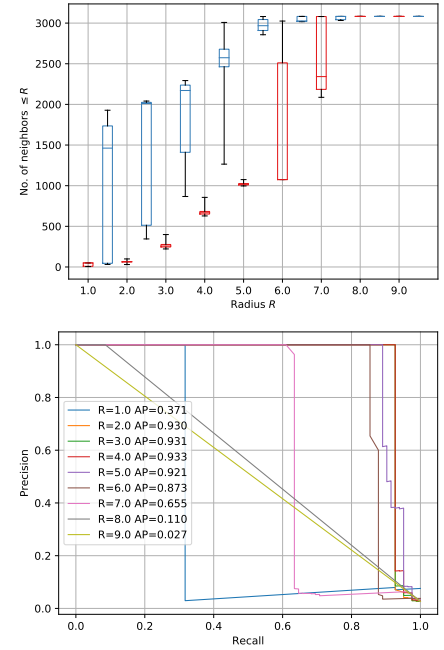


Figure 2: (Top) Number of neighbors within radius  $R$ , for anomalous (red) and benign (blue) points. (Bottom) PR curves using exact nearest neighbor counts.

### 1.3 Data-dependent Parameter Tuning

The anomaly detection performance is sensitive to picking the right parameters  $K$ ,  $T$  and  $w$  (equivalent to picking a radius  $R$  and the LSH approximation probabilities  $p_1$  and  $p_2$ ). There have been some approaches to auto-tuning the LSH parameters based on the data. However, they were primarily designed for nearest-neighbor search and do not work for data streams. Could we design a data structure for data streams that is tailored for anomaly detection?

- LSH-Forest (Bawa et al. 2005). The earliest approach; this is a very simple approach that builds a trie using each bit of the  $K$ -element sketch. Demonstrated for the Jaccard and cosine distances (since their LSH values are binary). A theoretical analysis was recently published (Andoni et al. 2017).
- Parameter-free LSH (Ahle et al. 2017). The most recent approach, designed for Euclidean distances. This builds a “multi-level LSH” data structure that tries multiple values of the projection size  $K$ , and finds the best level for each nearest-neighbor query.

### 1.4 Relationship to Chains

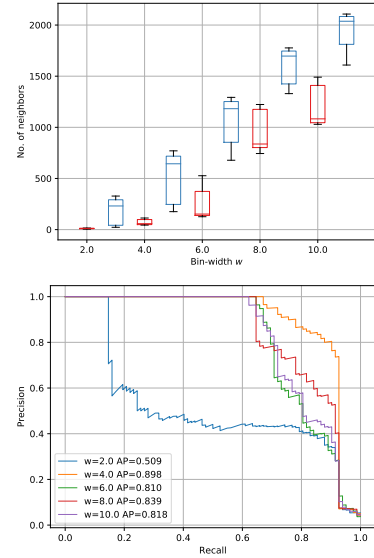


Figure 3: (Top) Distribution of approximate nearest neighbor counts via HF for different  $w$  (Bottom) PR curve for HF.  $K = 25$ ,  $T = 100$ .

## 2 Chains

**Pre-initialization.** This is done once, and is unchanged for all chains.

1. **Initialize the random projection matrix**  $A$ .
2. **Project the data.**  $Y = XA^T$ .
3. **Set**  $\Delta_f$  = half the range of projected feature  $f$ .
4. **Sample a shift**  $b_{if}$  for each chain  $i$  and feature  $f$ .

**Chain fitting.** For each chain  $i = 1, \dots, C$ :

For each depth  $d = 1, \dots, D$ :

1. **Sample a dimension**  $f_d \in \{1, \dots, K\}$  and store the dimension sampled at each depth. Let  $0 \leq c(f, d) \leq d$  be the number of times dimension  $f$  has been sampled in the chain until depth  $d$ .
2. **Bin the data.** Let  $Y_j = (y_{j1}, \dots, y_{jK})$  be a projected data point. Denote by the  $Z_{jd}[f]$  the “unfloored bin index” of feature  $f$  for point  $j$  at depth  $d$ :

$$Z_{jd}[f] = \frac{y_{jf} + b_{if}/2^{c(f,d)-1}}{\Delta_f/2^{c(f,d)-1}} \quad \forall f = 1, \dots, K \quad (1)$$

At depth  $d$ , the bin indices  $\tilde{Y}_{jd}$  are logically defined as follows, for  $f = 1, \dots, K$ :

$$\tilde{Y}_{jd}[f] = \begin{cases} 0, & \text{if } c(f, d) = 0 \\ \lfloor Z_{jd}[f] \rfloor & \text{if } c(f, d) > 0 \end{cases} \quad (2)$$

In practice, the bins are computed recursively at each depth as follows:

$$\tilde{Y}_{jd}[f] = \begin{cases} 0 & \text{if } c(f, d) = 0 \\ \lfloor Z_{jd}[f] \rfloor & \text{if } c(f, d) = 1 \\ \lfloor 2 \times Z_{jd-1}[f] - \frac{b_{if}}{\Delta_f} \rfloor & \text{if } c(f, d) > 1 \end{cases} \quad (3)$$

The recurrence is derived as follows:

$$x_k = \frac{y + b/2^k}{\Delta/2^k} = \frac{y + b/2^{k-1}}{\Delta/2^k} - \frac{b/2^k}{\Delta/2^k} = 2 \left( \frac{y + b/2^{k-1}}{\Delta/2^{k-1}} \right) - \frac{b}{\Delta} = 2x_{k-1} - \frac{b}{\Delta} \quad (4)$$

3. **Store bin counts.** In an exact hashtable or count-min sketch  $H_{id}$ , increment the count  $H_{id}(\tilde{Y}_j)$  for  $j = 1, \dots, N$ .

**Scoring at depth  $d$ .** The anomaly score of a point  $X_j$  at depth  $d$  is  $-\text{average}_i(H_{id}(\tilde{Y}_j))$ .

Chain Parameters	
$K$	Dimension of projected data
$C$	Number of chains
$D$	Chain depth
$\Delta_f$	Initial/maximum bin width
Other variables	
$f$	A feature index in $1, \dots, K$
$X$	$N \times D$ data matrix
$Y$	$N \times K$ projected data matrix
$A$	$K \times D$ projection matrix
$b_{if}$	Shift $\sim \text{Uniform}(0, \Delta_f)$
$H_{id}$	Hashtable/count-min sketch

Table 2: Notation.

## 2.1 Results with bin-counts as anomaly scores

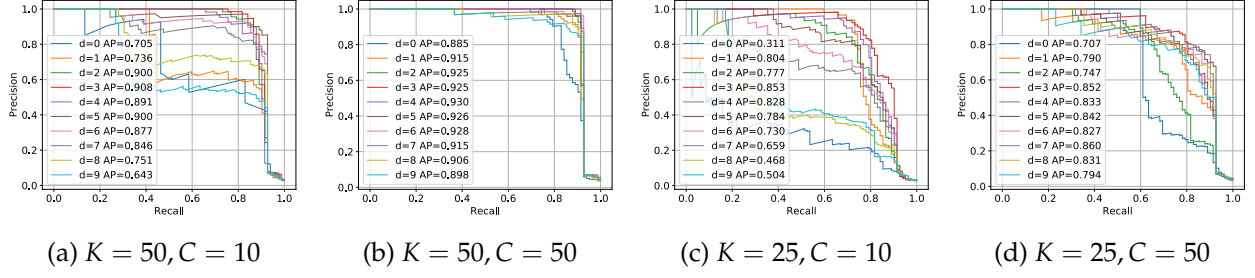


Figure 4: Precision-recall curves on synthetic data with 100 noisy dimensions,  $D = 10$ .

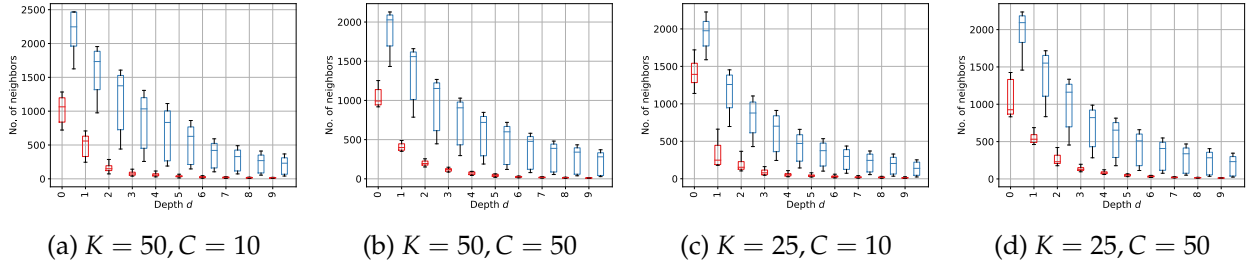


Figure 5: No. of neighbors in the same bin for anomalous (red) and benign (blue) points.

## 2.2 Results with LOCI scores as anomaly scores

Previously, we scored each point  $X_j$  using chain  $i$  at depth  $d$  using the number of neighbors lying in the same bin,  $H_{id}(\tilde{Y}_j)$ . We modify this to obtain a LOCI score  $L_{id}(\tilde{Y}_j)$  as:

$$\bar{H}_{id}(\tilde{Y}_j) = H_{id}(\tilde{Y}_j) \times 2^d \quad (5)$$

$$L_{id}(\tilde{Y}_j) = \bar{H}_{id}(\tilde{Y}_j) - \frac{1}{N-1} \sum_{k \neq j} \bar{H}_{id}(\tilde{Y}_k) \quad (6)$$

Then the anomaly score of a point  $X_j$  is  $-\text{average}_i(\min_d L_{id}(\tilde{Y}_j))$ .

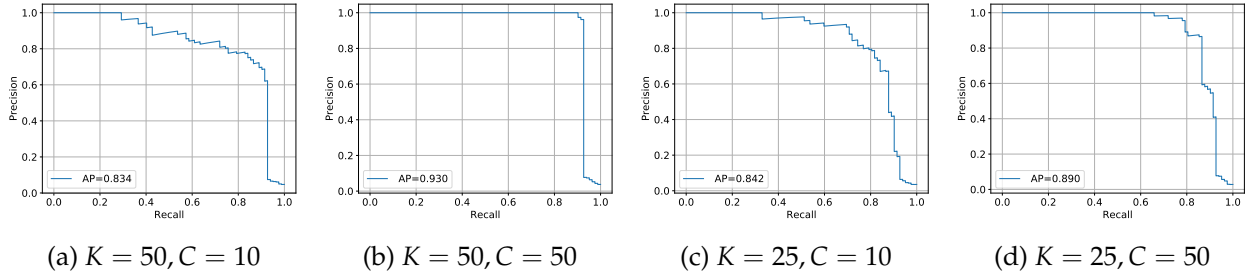


Figure 6: Precision-recall curves on synthetic data with 100 noisy dimensions,  $D = 10$ .

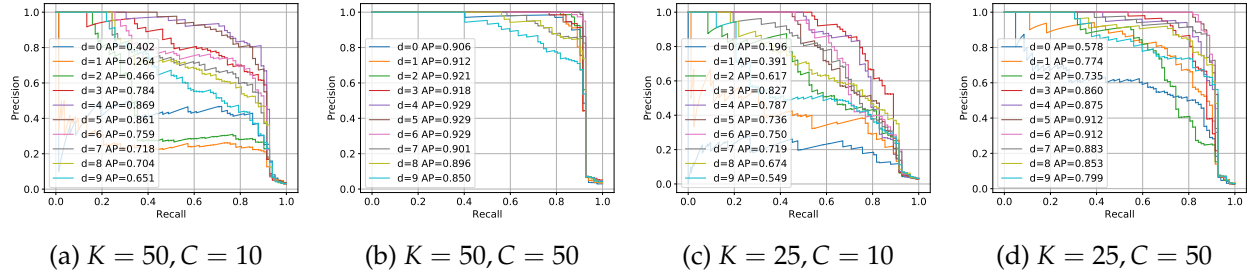


Figure 7: Precision-recall curves with LOCI scores at different depths  $d$ .

### 2.3 Stability across runs

Chain Parameters	Mean AP	Standard Deviation
$K = 50, C = 10$	0.824	0.064
$K = 50, C = 50$	0.930	0.001
$K = 25, C = 10$	0.742	0.123
$K = 25, C = 50$	0.901	0.019

Table 3: Chain AP stability across 5 runs,  $D = 10$ .

### 2.4 Number of unique bin IDs with depth

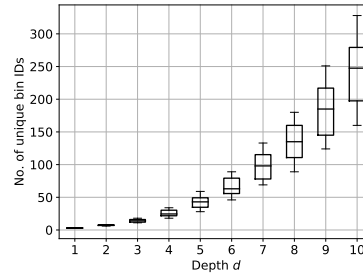


Figure 8: No. of unique bins for  $C = 100$  chains,  $K = 50, d = 10$ .

### 2.5 Scores for different anomaly types

Fig. 9 illustrates how the anomaly scores vary for different classes of points, shown in Fig. 9(a):

Class	No. of points	Description	Color
Class 0	1000	Benign sparse points	Blue
Class 1	2000	Benign dense points	Orange
Class 2	50	Anomalous dense points	Green
Class 3	25	Anomalous sparse points	Red
Class 4	6	Locally anomalous points	Purple
Class 5	1	Single anomalous point	Brown

From Fig. 9(b) conforms to intuition: as the depth in the chain increases, the average number of points lying in the same bin (across chains) decreases for all classes of points. This is due to two factors: (i) increasing the number of bin dimensions, and (ii) decreasing the bin width. Point classes lying in denser regions have higher bin counts.

Fig. 9(c) shows the LOCI scores (eq. 6) at each depth  $d$ , scaled up by  $2^d$  for visual clarity. The LOCI score of a point is simply its bin count minus the mean bin count of all other points at that depth. Hence, as  $d$  increases and the bin counts of all points tend to 1, the LOCI score tends to 1.

Fig. 9(d) shows the anomaly scores. The order of anomalousness across point classes is inversely proportional to their global density.

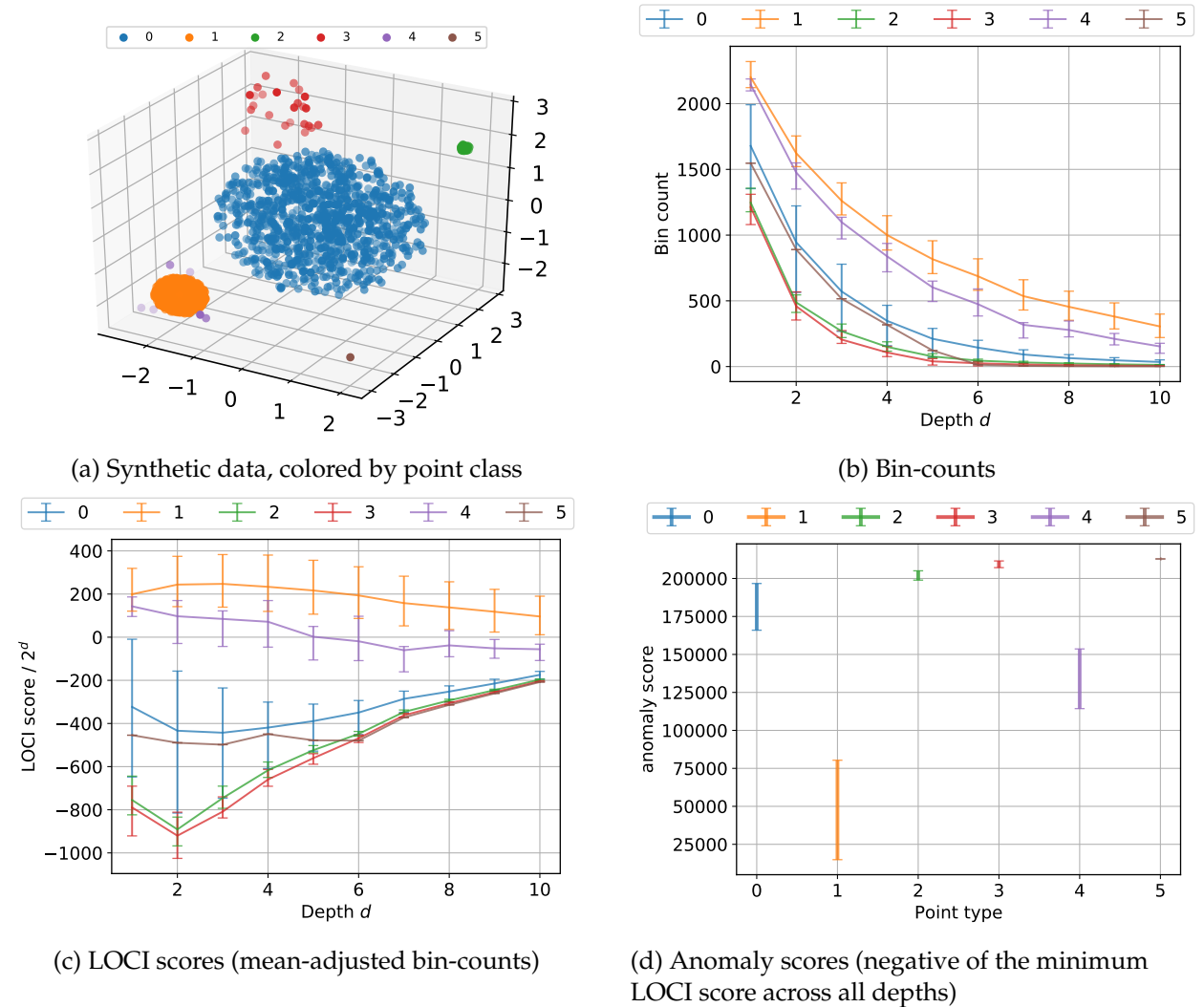


Figure 9: Score distributions for different point classes (see §2.3 for details). Error bars indicate the median, 5<sup>th</sup> and 95<sup>th</sup> percentile scores across all points. Chains with  $K = 50, C = 100, d = 10$ .

## 2.6 Streaming random projections

Let the data dimensionality be  $D$  and projection dimensionality be  $K$ . Given a density  $s$ , the very sparse random projections method (Li et al. 2006) generates  $K$   $D$ -dimensional vectors containing elements  $r_{km}$  as follows ( $k = 1, \dots, K, m = 1, \dots, D$ ):

$$r_{km} = \begin{cases} -\frac{1}{\sqrt{Ks}} & \text{with probability } \frac{s}{2} \\ 0 & \text{with probability } 1 - s \\ \frac{1}{\sqrt{Ks}} & \text{with probability } \frac{s}{2} \end{cases}$$

Instead of storing these vectors, we store  $K$  hash functions. Each hash function  $h_k$  hashes a string  $w_m$  (the feature name) to a hash value  $h_k(w_m)$ , which serves as a plug-in replacement for  $r_{km}$  without explicitly storing the random vectors.

The implementation of  $h_k$  ensures that the hash values follow the probability distribution prescribed above.

The function  $h_k(w_m)$  is implemented as follows:

1.  $b = g_k(w_m)$ , where  $g_k$  is drawn from a universal family and  $b$  is a 32-bit integer. Having more bits in  $b$  provides a more precise approximation to the desired probability distribution, but we found 32 to be sufficient.
2.  $n = b / (2^{32} - 1)$  is a number between 0 and 1. With a good universal family to generate  $g_k$ ,  $n$  follows distributions  $P_{g_k}(n) \sim \text{Uniform}(0, 1)$  and  $P_{w_m}(n) \sim \text{Uniform}(0, 1)$ .
3.  $n$  is used to return a number in  $\{\frac{1}{\sqrt{Ks}}, 0, \frac{1}{\sqrt{Ks}}\}$  with the desired probability distribution:
  - (a) If  $0 \leq n < s/2$ , return  $-\frac{1}{\sqrt{Ks}}$ .
  - (b) If  $s/2 \leq n < s$ , return  $\frac{1}{\sqrt{Ks}}$ .
  - (c) If  $s \leq n \leq 1$ , return 0.

It is very important that  $g_k$  be drawn from a universal (or almost-universal) string hash family with good empirical randomization properties. SipHash (Aumasson et al. 2012) and MurmurHash3<sup>1</sup> were found to perform well, whereas multiplicative hashing does not.

Fig. 10 illustrates the empirical distortion of this streaming random projection method on the synthetic data, for both the Euclidean and cosine distances. Fig. 11 illustrates the distortion using the standard, static random projections.

When used as a plug-in replacement for sparse random projections in XStream on the synthetic data, there was no significant drop in average precision (mean and standard deviation reported for 5 runs): AP =  $0.823 \pm 0.088$  for  $K = 50, C = 50$ ; AP =  $0.916 \pm 0.009$  for  $K = 100, C = 50$ ; AP =  $0.927 \pm 0.003$  for  $K = 100, C = 100$ .

<sup>1</sup><https://github.com/aappleby/smhasher/blob/master/src/MurmurHash3.cpp>



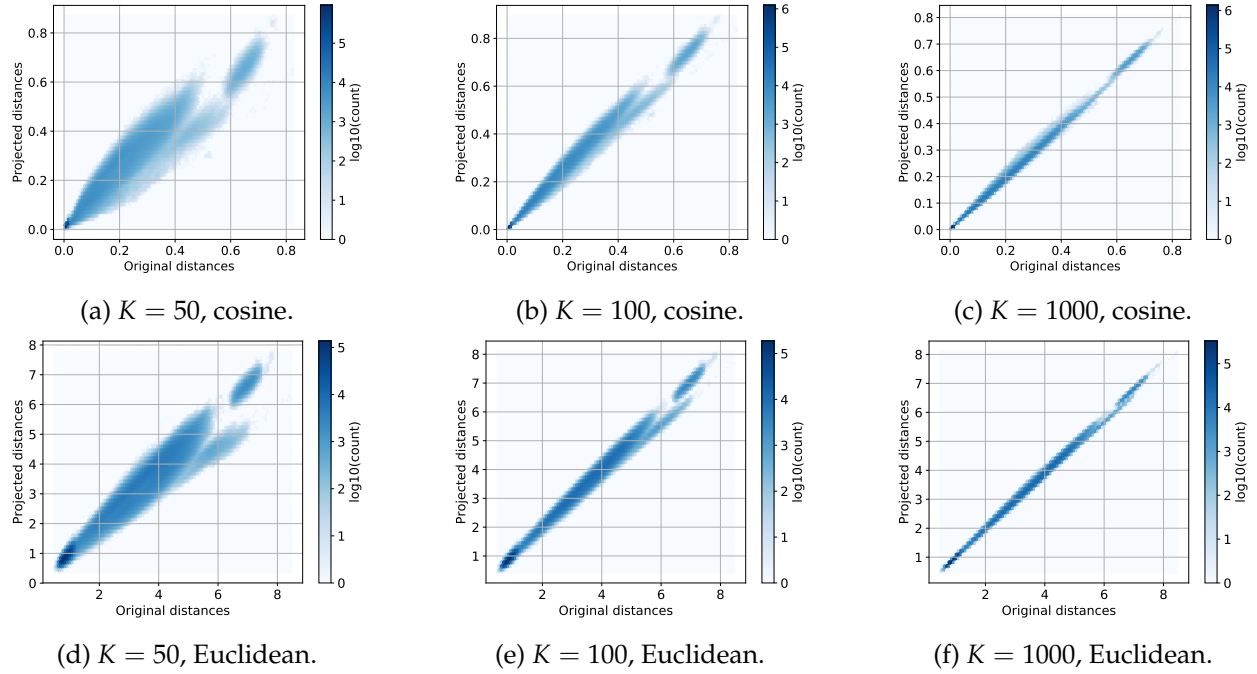


Figure 10: Distortion of streaming random projections for the Euclidean and cosine distance.

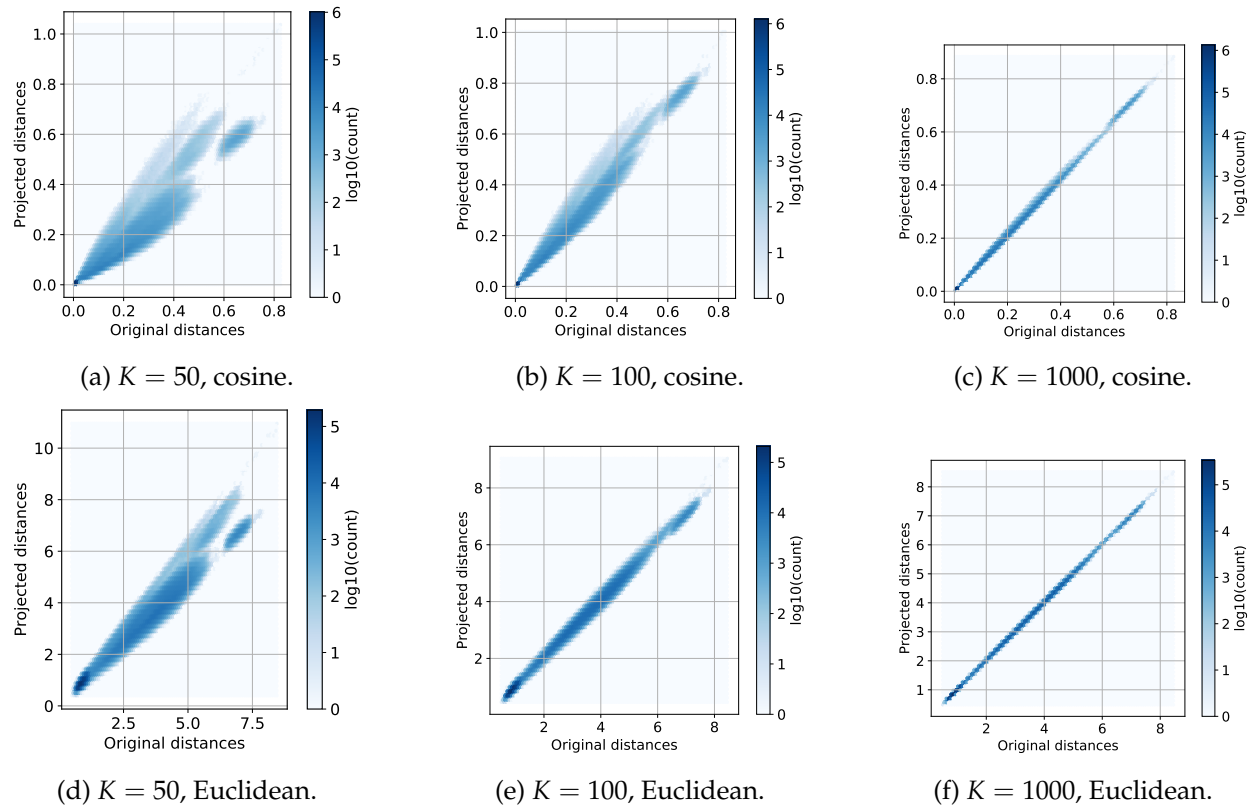


Figure 11: Distortion for static random projections for the Euclidean and cosine distance.

## 2.7 Addressing local density variation

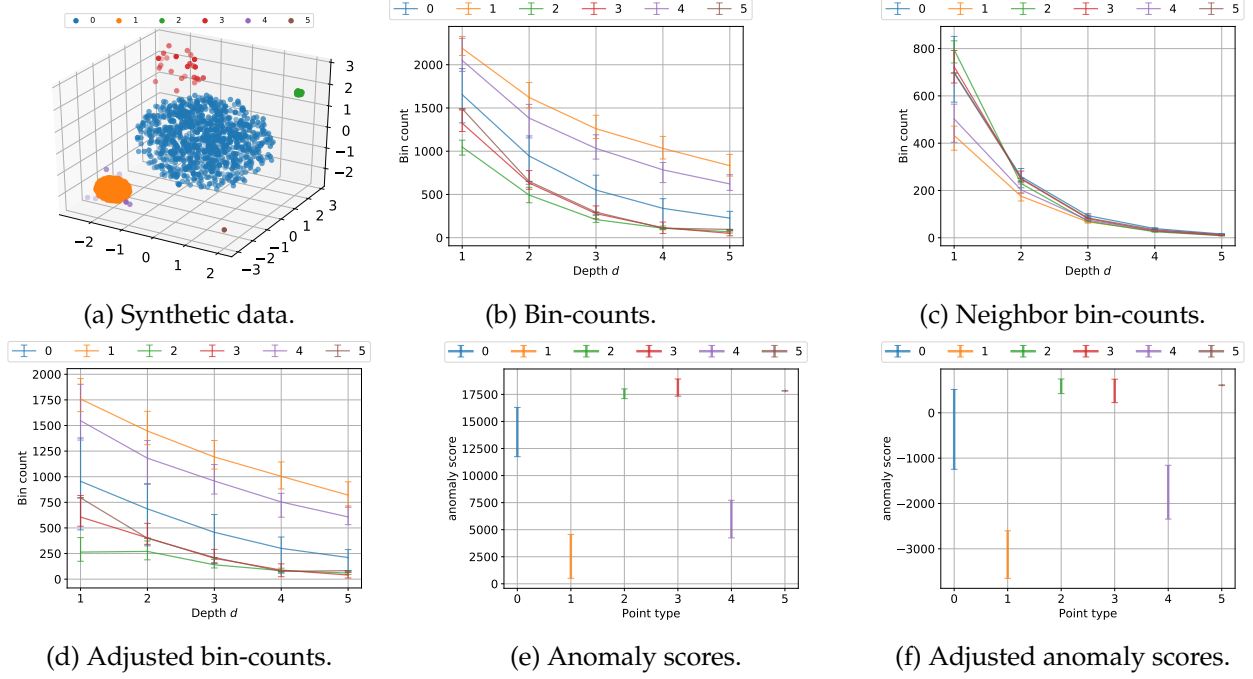


Figure 12: Score distributions for different point classes after addressing local density variation (see §2.3 and §2.7). Error bars indicate the median, 5<sup>th</sup> and 95<sup>th</sup> percentile scores across all points. Chains with  $K = 50, C = 100, d = 5$ .

Instead of adjusting the bin-count of each point  $j$  with the mean bin-count of all other points, we could adjust it with the mean bin-count over *neighboring bins*. This enables us to capture deviations in the *local* density.

However, for a given point  $j$  having a  $d$ -element bin-id at depth  $d$ , there are  $3^d - 1$  possible neighboring bins; enumerating all of them is infeasible, and most of them would be empty.

Instead, we could sample a fixed number of neighboring bins at random, such that the probability of sampling each neighboring bin is proportional to the likelihood of neighbors of  $j$  lying in that bin. This probability is proportional to the distance of the point  $j$  to the edge of the neighboring bin (see  $l_1$  and  $l_2$  in Fig. 13: shaded neighboring bins are more likely to contain neighbors of the red point  $j$ ). This idea was introduced as “multi-probe LSH” in (lv2007multi).

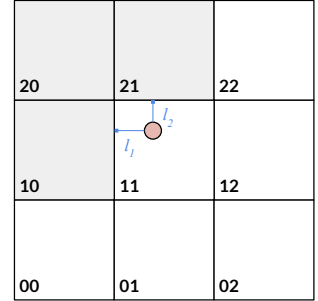


Figure 13: 2-d grid.

Thus, for a point  $j$  having a set of neighboring bins  $\mathcal{N}(j)$ , the adjusted LOCI score becomes:

$$L_{id}(\tilde{Y}_j) = \bar{H}_{id}(\tilde{Y}_j) - \frac{1}{|\mathcal{N}(j)|} \sum_{\tilde{Y}_k \in \mathcal{N}(j)} \bar{H}_{id}(\tilde{Y}_k) \quad (7)$$

## 2.8 Incremental maintenance of the LOCI score

The LOCI score for chain  $i$ , point  $j$  and depth  $d$  was defined as:

$$L_{id}(\bar{Y}_j) = H_{id}(\bar{Y}_j) \times 2^d - \frac{1}{N} \sum_k H_{id}(\bar{Y}_k) \times 2^d \quad (8)$$

where  $H_{id}(\bar{Y}_j)$  is the bin-count and is directly available on the arrival of every new point. Hence, the first term of the above equation is easily computable.

Let  $b_{id}(c)$  be the bin-count of a bin  $c$  in chain  $i$  at depth  $d$ . Then the total bin-count  $S_{id}$  (the sum of the bin-counts of all points in chain  $i$  at depth  $d$ ) is given by:

$$S_{id} = \sum_k H_{id}(\bar{Y}_k) = \sum_{c \in \text{bins in chain}(i,d)} b_{id}(c)^2 \quad (9)$$

When a new point arrives, it falls into some cell  $c'$ , and its bin-count is incremented to  $b_{id}(c') + 1$ . This results in the following change to the total bin-count:

$$S'_{id} = \sum_{c \neq c'} b_{id}(c)^2 + (b_{id}(c') + 1)^2 \quad (10)$$

$$= \sum_c b_{id}(c)^2 + 2b_{id}(c') + 1 \quad (11)$$

$$= S_{id} + 2b_{id}(c') + 1 \quad (12)$$

where  $b_{id}(c)$  was the bin-count before the point was added, and is easily computable. This equation, along with tracking the number of points  $N$ , can thus be used to incrementally maintain the second term of the LOCI score.

## 2.9 Dynamic chains with windows

## 2.10 Streaming evaluation

Dataset	Feature space	Dynamic updates	Input size (incl. updates)	Dimensionality (maximum)	Anomaly rate
network-intrusion	Fixed	No	720K	34	00.87%
spam-sms	Evolving	No	5574	8442	13.40%
spam-url	Evolving	No	3.2M	2.4M	33.00%
streamspot-flash	Evolving	Yes	111M	?	20.00%
streamspot-java	Evolving	Yes	?	?	20.00%

Table 4: Datasets used for the streaming evaluation.

### Baselines

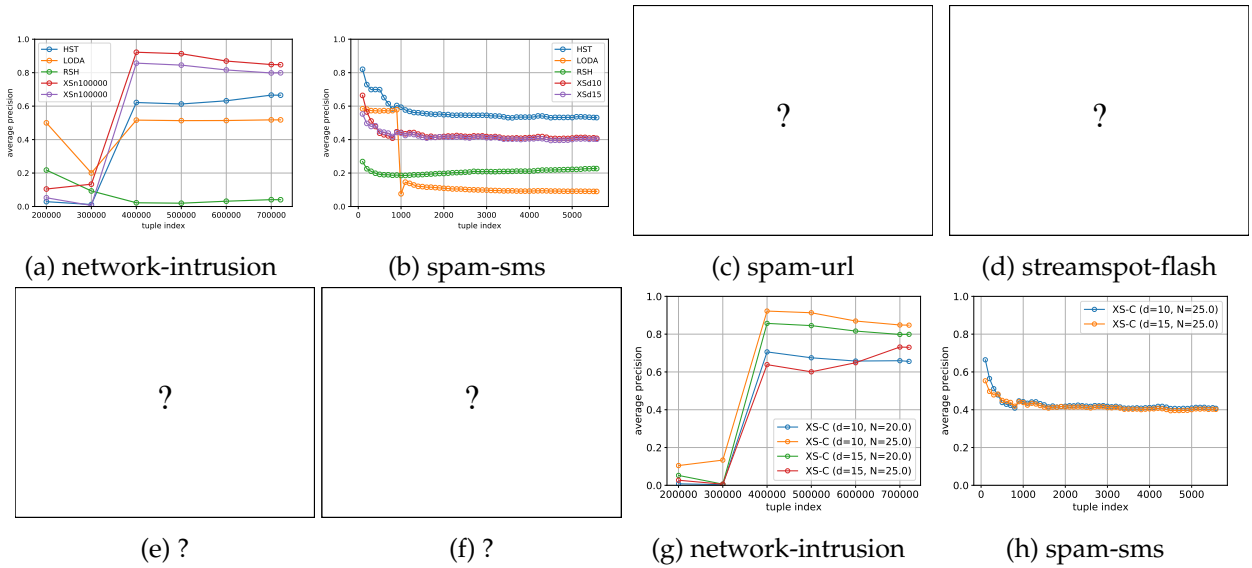


Figure 14: (a)-(e) Average precision over time for XS-C, XS-W and baseline methods for 5 datasets. (f)-(h) Effect of varying  $d$  and initial sample size  $N$  on XS-C.

### Observations.

1. There may be some sensitivity to the depth  $d$  for XS-C. This reduces as  $C$  is increased.
2. Baseline method implementations have been verified to give the right ROC-AUC numbers on the datasets they originally reported on.

Methods	Window / Continuous	Parameters
HS-Trees	Window	Max Depth=15, Num Trees=100, Window Size = 256 Decay Rate = 0.015, Initial Sample Size = 25% Two-Window Histogram, WindowSize = 256
RS-Hash	Continuous	
LODA	Window	

Table 5: Baseline Methods and the parameters chosen based on best performance

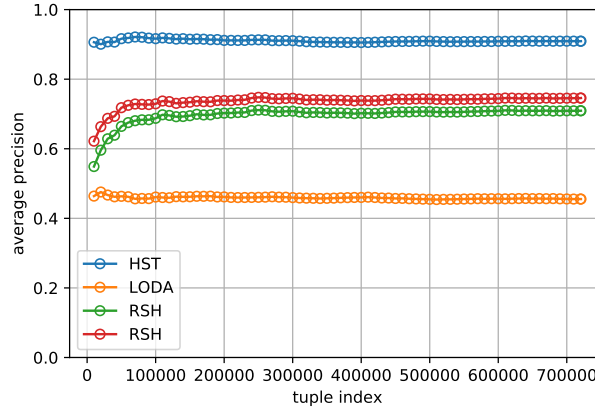


Figure 15: RS-Hash performance improves on shuffled version of HTTP dataset.

3. High ROC-AUC as reported by baseline methods is not correlated with high AP (see <http://www.dbs.ifi.lmu.de/research/outlier-evaluation/DAMI/literature/KDDCup99/> for an example on the network-intrusion dataset). In fact, AP is a more sensitive (and better) metric than ROC-AUC.
4. RSH performs poorly on network-intrusion due to clustered anomalies. On shuffling the dataset, the accuracy improves ??.

**PR-AUC (AP) vs. ROC-AUC.** PR curves do not care about “true negatives”, which do not show up in the formulae for precision or recall. Hence, there is no advantage in labelling non-anomalies correctly: this is desired, since having a large number of non-anomalies (as is usually the case) reduces the loss in AUC due to incorrectly labelling anomalies. In fact, a curve that dominates in PR-space will also dominate in ROC-space, but the converse is not true (**davis2006relationship**).

### 3 Static evaluation

Dataset	Number of samples	Dimensionality	Number of Anomalies
<i>High-dimensional Datasets</i>			
gisette	7000	4970	3500
isolet	7797	616	3300
letter	7797	616	3600
madelon	2600	500	1300
<i>Low/medium-dimensional Datasets</i>			
cancer	568	30	211
ionosphere	350	33	125
telescope	19020	10	6688
indians	768	8	269

Table 6: Datasets used for the static evaluation.

### 3.1 Datasets

The datasets used are listed in Table 6 to generate the subsequent datasets used in the static evaluation. They were chosen based on (TODO: selection logic).

### 3.2 Generating anomalies for static evaluation

We used 2 different procedures to generate the datasets on which to benchmark, one for the low/medium dimensional datasets and one for the high-dimensional datasets. They are described below:

1. *Approach 1:* We append new columns containing Gaussian noise to the original dataset, thereby increasing its dimensionality with noisy features. The number of noisy features is quantified as a percentage of the original dataset dimensionality. If the original dataset mean and standard deviation are  $\mu$  and  $\sigma$ , the noisy columns contained Gaussian noise with mean  $\rho\mu$  and standard deviation  $\rho\sigma$ , where  $\rho$  is the noise factor. We consider appending 100%, 1000%, 2000% and 5000% noisy features, and use  $\rho \in \{0.01, 0.1, 0.2, 0.25\}$ .
2. *Approach 2:* We first discard all the original anomalous points. Now for the remaining nominal points, we first choose 10% of the features at random and mark them as important features. We further proceed with choosing 5% of the samples at random as important samples. Now, in this sub-block, which contains only important samples and important features, we add Gaussian noise to the original data. Again, Gaussian noise is used, parameterized by the mean, and the standard deviation of the entire nominal dataset. Further this noise is scaled according to the user-specified signal to noise ratio, and added to the original sub-block of our dataset and features.

### 3.3 Competing methods and hyperparameter settings

**XStream (XS)** was used with  $K = 100$  projections,  $C = 100$  chains and depth  $d = 15$ . We implemented and compared against the following competing methods and parameter settings:

- **iForest (IF)**: We use it with recommended sample size of 256 or the entire dataset, whichever is minimum. The `hlim` is set to 15 to be comparable to HS-Trees, and number of components is set to 100.
- **HS-Trees (HST)**: Use it with the recommended setting of 15 max-depth and 100 trees.
- **LODA**: We set the sparsity factor to  $\frac{1}{\sqrt{(d)}}$ . We allow LODA to select best width for histograms, but fix the number of histograms to 100.
- **RS-Hash (RSH)**: We use 1000 sampling points and 100 components.

### 3.4 Baseline Results

Table 9 lists the average precision of XStream and the competing methods on the generated datasets.

We conduct our experiments on datasets mentioned in Table 6. A Friedman test for differences in the best-performing method across all datasets rejected the null hypothesis that the difference rankings between methods is not statistically significant with  $p = 0.0072 < 0.05$ .

Hence, we perform a posthoc-Friedman test, Nemenyi test to compare all methods to each other. We first compute average ranks of all the methods over the 8 original datasets, which is shown in Table 7. Setting significance level to be  $\alpha = 0.05$ , we get the  $q_{0.05}$  for  $k = 5$  pairs as 2.728. Difference in average ranks of any two methods will be significant if

$$(R_i - R_j) > q_\alpha \sqrt{\frac{k(k+1)}{6N}} \quad (13)$$

Setting  $N = 8$  and  $k = 5$ , for difference between a pair of methods to be significant, the difference between the rank should be greater than 2.1567. Looking at the average ranks shown in Table 7, we can say the (X-Stream, LODA) and (X-Stream, I-Forest) are significantly different.

To show the effect of adding high-dimensionality, we add noisy columns to the 4 original datasets. The noise is generated as mentioned above in Approach 1. The average ranks for all the methods is shown in Table 8. Using Friedman test, we can reject the null hypothesis that the difference rankings between methods is not statistically significant with  $p = 3.451 \times 10^{-10}$ . We again perform posthoc-Friedman Nemenyi test; setting  $N = 16$  and  $k = 5$ , we obtain that difference in average ranks should be greater than 1.524. Observing the table, we find difference to be significant between (LODA, I-Forest), (X-Stream, I-Forest), (LODA, HS-Trees), (X-Stream, HS-Trees) and (X-Stream, RS-Hash).

	IF	HST	RSH	LODA	XS
cancer	0.635(5.0)	0.758(3.0)	0.719(4.0)	0.883(2.0)	0.9035(1.0)
ionosphere	0.819(4.0)	0.824(3.0)	0.83(2.0)	0.79(5.0)	0.8952(1.0)
telescope	0.655(3.0)	0.674(2.0)	0.678(1.0)	0.623(5.0)	0.6339(4.0)
indians	0.498(4.0)	0.502(3.0)	0.518(2.0)	0.48(5.0)	0.5382(1.0)
gisette	0.423(4.0)	0.417(5.0)	0.432(3.0)	0.436(2.0)	0.4424(1.0)
isolet	0.442(4.0)	0.444(2.5)	0.444(2.5)	0.429(5.0)	0.4445(1.0)
letter	0.475(2.0)	0.453(5.0)	0.471(3.0)	0.46(4.0)	0.4865(1.0)
madelon	0.501(5.0)	0.507(4.0)	0.511(3.0)	0.515(2.0)	0.5277(1.0)
AVG RANK	3.875	3.4375	2.5625	3.75	1.375

Table 7: Average rank of methods over the 8 original datasets

	IF	HST	RSH	LODA	XS
cancer100	0.661(5.0)	0.727(3.0)	0.702(4.0)	0.839(2.0)	0.8511(1.0)
cancer1000	0.476(4.0)	0.394(5.0)	0.531(3.0)	0.686(2.0)	0.8453(1.0)
cancer2000	0.447(4.0)	0.444(5.0)	0.474(3.0)	0.643(2.0)	0.8282(1.0)
cancer5000	0.404(4.0)	0.4(5.0)	0.425(3.0)	0.508(2.0)	0.7729(1.0)
ionosphere100	0.756(4.0)	0.773(2.0)	0.725(5.0)	0.771(3.0)	0.9005(1.0)
ionosphere1000	0.59(3.0)	0.573(5.0)	0.579(4.0)	0.742(2.0)	0.8706(1.0)
ionosphere2000	0.441(5.0)	0.505(3.0)	0.446(4.0)	0.736(2.0)	0.8743(1.0)
ionosphere5000	0.403(3.0)	0.376(5.0)	0.399(4.0)	0.675(2.0)	0.7966(1.0)
telescope100	0.593(4.0)	0.576(5.0)	0.601(3.0)	0.617(2.0)	0.6373(1.0)
telescope1000	0.451(4.0)	0.419(5.0)	0.471(3.0)	0.592(2.0)	0.6095(1.0)
telescope2000	0.407(4.5)	0.407(4.5)	0.409(3.0)	0.586(2.0)	0.5872(1.0)
telescope5000	0.376(5.0)	0.38(3.0)	0.378(4.0)	0.538(2.0)	0.5752(1.0)
indians100	0.477(3.0)	0.444(5.0)	0.478(2.0)	0.472(4.0)	0.5218(1.0)
indians1000	0.385(4.0)	0.379(5.0)	0.397(3.0)	0.444(2.0)	0.4976(1.0)
indians2000	0.35(5.0)	0.361(4.0)	0.369(3.0)	0.41(2.0)	0.4777(1.0)
indians5000	0.343(4.0)	0.338(5.0)	0.349(3.0)	0.4(2.0)	0.4659(1.0)
AVG. Rank	4.09375	4.3475	3.375	2.1875	1

Table 8: Average rank of methods over the 16 perturbed datasets.



Dataset	IF	HST	RSH	LODA	XS
<i>High-dimensional Datasets</i>					
gisette(30,0.3,1.2 )	$0.683 \pm 0.045$	$0.46 \pm 0.005$	$0.628 \pm 0.021$	$0.531 \pm 0.018$	$0.5282 \pm 0.0091$
gisette (30,0.3,10.0 )	$0.541 \pm 0.027$	$0.354 \pm 0.002$	$0.471 \pm 0.024$	$0.321 \pm 0.004$	$0.3204 \pm 0.0029$
gisette (30,0.3,20.0 )	$0.45 \pm 0.022$	$0.308 \pm 0.003$	$0.347 \pm 0.017$	$0.29 \pm 0.003$	$0.292 \pm 0.0032$
gisette (30,0.3,30.0 )	$0.432 \pm 0.016$	$0.31 \pm 0.002$	$0.316 \pm 0.006$	$0.301 \pm 0.005$	$0.2929 \pm 0.0023$
isolet	$0.442 \pm 0.019$	$0.444 \pm 0.001$	$0.444 \pm 0.013$	$0.429 \pm 0.02$	$0.4445 \pm 0.0067$
isolet (30,0.3,1.2 )	$0.537 \pm 0.03$	$0.441 \pm 0.01$	$0.559 \pm 0.028$	$0.58 \pm 0.03$	$0.6396 \pm 0.0204$
isolet (30,0.3,10.0 )	$0.372 \pm 0.011$	$0.334 \pm 0.004$	$0.391 \pm 0.015$	$0.35 \pm 0.007$	$0.3479 \pm 0.0066$
isolet (30,0.3,20.0 )	$0.33 \pm 0.009$	$0.311 \pm 0.001$	$0.322 \pm 0.007$	$0.313 \pm 0.006$	$0.3176 \pm 0.0017$
isolet (30,0.3,30.0 )	$0.3 \pm 0.002$	$0.296 \pm 0.001$	$0.293 \pm 0.002$	$0.292 \pm 0.004$	$0.2899 \pm 0.0016$
letter (30,0.3,1.2 )	$0.49 \pm 0.026$	$0.44 \pm 0.01$	$0.519 \pm 0.053$	$0.534 \pm 0.03$	$0.5967 \pm 0.0204$
letter (30,0.3,10.0 )	$0.374 \pm 0.015$	$0.339 \pm 0.008$	$0.375 \pm 0.007$	$0.337 \pm 0.003$	$0.3361 \pm 0.005$
letter (30,0.3,20.0 )	$0.309 \pm 0.008$	$0.291 \pm 0.002$	$0.31 \pm 0.007$	$0.289 \pm 0.002$	$0.3 \pm 0.0015$
letter (30,0.3,30.0 )	$0.319 \pm 0.01$	$0.308 \pm 0.001$	$0.307 \pm 0.004$	$0.305 \pm 0.003$	$0.3095 \pm 0.0023$
madelon (5,0.05,1.2 )	$0.896 \pm 0.051$	$0.92 \pm 0.007$	$0.969 \pm 0.039$	$1.0 \pm 0.0$	$1.0 \pm 0.0$
madelon (5,0.05,10.0 )	$0.643 \pm 0.13$	$0.784 \pm 0.068$	$0.899 \pm 0.094$	$0.988 \pm 0.012$	$0.9567 \pm 0.023$
madelon (5,0.05,20.0 )	$0.24 \pm 0.093$	$0.16 \pm 0.039$	$0.377 \pm 0.143$	$0.112 \pm 0.027$	$0.125 \pm 0.0133$
madelon (5,0.05,30.0 )	$0.07 \pm 0.013$	$0.083 \pm 0.004$	$0.099 \pm 0.059$	$0.054 \pm 0.007$	$0.0444 \pm 0.0034$
<i>Low/medium-dimensional Datasets</i>					
cancer (100,0.1 )	$0.661 \pm 0.036$	$0.727 \pm 0.01$	$0.702 \pm 0.026$	$0.839 \pm 0.024$	$0.8511 \pm 0.0216$
cancer (1000,0.1 )	$0.476 \pm 0.043$	$0.394 \pm 0.007$	$0.531 \pm 0.047$	$0.686 \pm 0.047$	$0.8453 \pm 0.0179$
cancer (2000,0.1 )	$0.447 \pm 0.026$	$0.444 \pm 0.015$	$0.474 \pm 0.046$	$0.643 \pm 0.066$	$0.8282 \pm 0.0212$
cancer (5000,0.1 )	$0.404 \pm 0.021$	$0.4 \pm 0.013$	$0.425 \pm 0.022$	$0.508 \pm 0.086$	$0.7729 \pm 0.0394$
ionosphere (100,0.1 )	$0.756 \pm 0.031$	$0.773 \pm 0.005$	$0.725 \pm 0.018$	$0.771 \pm 0.013$	$0.9005 \pm 0.0052$
ionosphere (1000,0.1 )	$0.59 \pm 0.044$	$0.573 \pm 0.008$	$0.579 \pm 0.054$	$0.742 \pm 0.032$	$0.8706 \pm 0.01$
ionosphere (2000,0.1 )	$0.441 \pm 0.033$	$0.505 \pm 0.033$	$0.446 \pm 0.027$	$0.736 \pm 0.029$	$0.8743 \pm 0.0052$
ionosphere (5000,0.1 )	$0.403 \pm 0.022$	$0.376 \pm 0.006$	$0.399 \pm 0.028$	$0.675 \pm 0.048$	$0.7966 \pm 0.0117$
telescope (100,0.1 )	$0.593 \pm 0.01$	$0.576 \pm 0.005$	$0.601 \pm 0.015$	$0.617 \pm 0.003$	$0.6373 \pm 0.0077$
telescope (1000,0.1 )	$0.451 \pm 0.018$	$0.419 \pm 0.008$	$0.471 \pm 0.021$	$0.592 \pm 0.013$	$0.6095 \pm 0.0054$
telescope (2000,0.1 )	$0.407 \pm 0.021$	$0.407 \pm 0.003$	$0.409 \pm 0.018$	$0.586 \pm 0.007$	$0.5872 \pm 0.0087$
telescope (5000,0.1 )	$0.376 \pm 0.011$	$0.38 \pm 0.002$	$0.378 \pm 0.005$	$0.538 \pm 0.026$	$0.5752 \pm 0.0074$
indians (100,0.1 )	$0.477 \pm 0.013$	$0.444 \pm 0.009$	$0.478 \pm 0.007$	$0.472 \pm 0.01$	$0.5218 \pm 0.0077$
indians (1000,0.1 )	$0.385 \pm 0.017$	$0.379 \pm 0.006$	$0.397 \pm 0.024$	$0.444 \pm 0.014$	$0.4976 \pm 0.0106$
indians (2000,0.1 )	$0.35 \pm 0.014$	$0.361 \pm 0.005$	$0.369 \pm 0.023$	$0.41 \pm 0.016$	$0.4777 \pm 0.0082$
indians (5000,0.1 )	$0.343 \pm 0.008$	$0.338 \pm 0.003$	$0.349 \pm 0.014$	$0.4 \pm 0.028$	$0.4659 \pm 0.009$

Table 9: Average precision of static methods on high-dimensional (top) and low/medium-dimensional (bottom) datasets. Mean and standard deviation reported over 10 runs. Numbers in the brackets indicate: (top) the percentage of features, fraction of samples to which noise is added, signal-to-noise ratio, and (bottom) noise column amount (as % of original dimensionality), relative noise factor.

Dataset	IF	HST	RSH	LODA	XS
<i>Original Datasets</i>					
gisette	$0.423 \pm 0.011$	$0.417 \pm 0.005$	$0.432 \pm 0.017$	$0.436 \pm 0.008$	$0.4424 \pm 0.0081$
isolet	$0.442 \pm 0.019$	$0.444 \pm 0.001$	$0.444 \pm 0.013$	$0.429 \pm 0.02$	$0.4445 \pm 0.0067$
letter	$0.475 \pm 0.012$	$0.453 \pm 0.003$	$0.471 \pm 0.009$	$0.46 \pm 0.011$	$0.4865 \pm 0.0075$
madelon	$0.501 \pm 0.008$	$0.507 \pm 0.003$	$0.511 \pm 0.01$	$0.515 \pm 0.014$	$0.5277 \pm 0.008$
cancer	$0.635 \pm 0.024$	$0.758 \pm 0.004$	$0.719 \pm 0.019$	$0.883 \pm 0.019$	$0.9035 \pm 0.0104$
ionosphere	$0.819 \pm 0.008$	$0.824 \pm 0.003$	$0.83 \pm 0.005$	$0.79 \pm 0.014$	$0.8952 \pm 0.0084$
telescope	$0.655 \pm 0.011$	$0.674 \pm 0.003$	$0.678 \pm 0.005$	$0.623 \pm 0.004$	$0.6339 \pm 0.0061$
indians	$0.498 \pm 0.015$	$0.502 \pm 0.005$	$0.518 \pm 0.006$	$0.48 \pm 0.015$	$0.5382 \pm 0.0052$

Table 10: Average precision of static methods on original, unperturbed datasets. Mean and standard deviation are reported over 10 runs.

Table 11: Symbol List

Symbols	Description
$N$	Number of data points
$D$	Number of dimensions/features
$C$	Number of ensemble components (trees, chains, histograms, etc.)
$d$	Max depth of trees/chains
$\psi$	iForest/HS-Trees/RS-Hash: sampling size
$r$	RS-Hash: Number of sampled dimensions
$b$	LODA: number of histogram bins
$m$	X-STREAM/RS-Hash: Number of CMS hash functions
$L$	X-STREAM/RS-Hash:: CMS hash table size
$k$	X-STREAM: Number of random projections
$M$	X-STREAM: Number of subspaces
$W$	X-STREAM: Number of windows

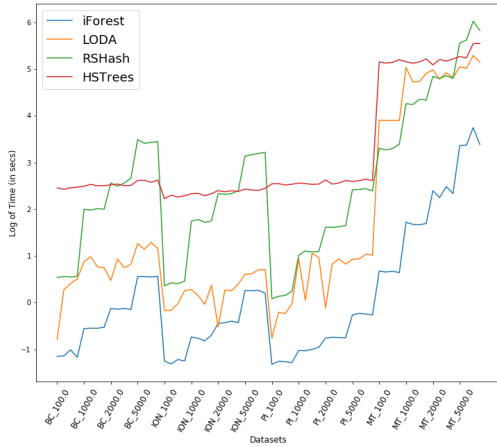
### 3.5 Time and Space Complexity

### 3.6 Running time

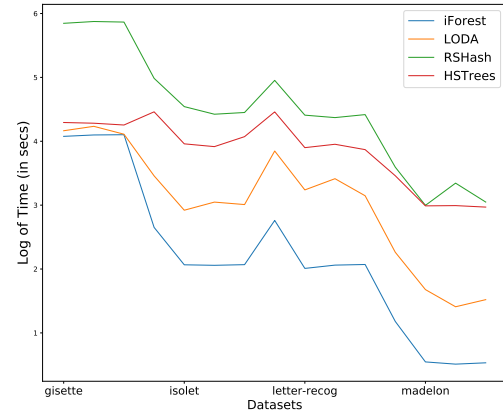
iForest is the fastest, RSHash is the slowest (RSHash implementation, however can be optimized)

Table 12: Time and space complexity of the ensemble anomaly detection algorithms compared in this paper. In streaming case, a data point/vector arrives at a time for HS-Stream, LODA, and RS-Hash; whereas an update to a single feature for a data point arrives at a time.

Algorithm	Time Complexity		Space Complexity
	Training	Scoring/Updating	
<i>Batch/Offline</i>			
iForest	$O(C\psi \log \psi)$	$O(C \log \psi)$	$O(C\psi)$
HS-Trees	$O(Cd\psi)$	$O(Cd)$	$O(C2^d)$
LODA	$O(NC\sqrt{D})$	$O(C\sqrt{D})$	$O(C\sqrt{D} + Cb)$
RS-Hash	$O(Crm\psi)$	$O(Crm)$	$O(CLm)$
X-STREAM	$O(M[NDk + Ckmd])$	$O(M[Dk + Ckmd])$	$O(MCLmd)$
<i>Streaming/Online</i>			
HS-Stream	—	$O(Cd + C\psi)$	$O(C2^d)$
LODA	—	$O(C\sqrt{D} + Cb)$	$O(C\sqrt{D} + Cb)$
RS-Hash	—	$O(Crm)$	$O(CLm)$
X-STREAM	—	$O(WMCKmd)$	$O(WMCLmd)$



(a) Low /Mid - dimensional dataset analysis



(b) High dimensional dataset analysis

Figure 16: RunTimes of algorithms over datasets with different noise levels. X-axis is over different datasets with different noise levels. Y-axis is log of runtime (in seconds).

## 4 Dynamic evaluation

### Datasets:

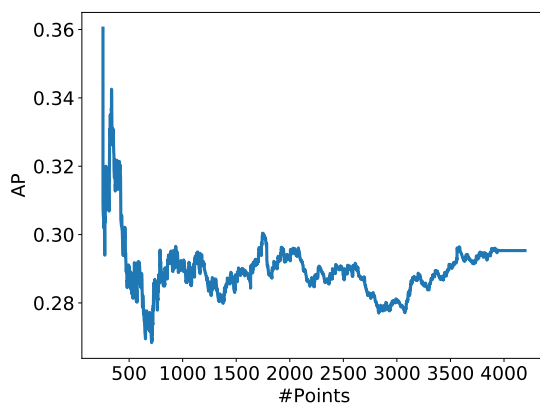
We use the same datasets as in HighDimension case in static setting. We create 10 shuffled versions for each dataset. The datasets could either be shuffled randomly or it could be clustered in time. For the clustered shuffled, the input also includes a clustering percentage, that is how much percentage of anomalies should be clustered in time. We currently do not experiment with

clustered in time datasets.

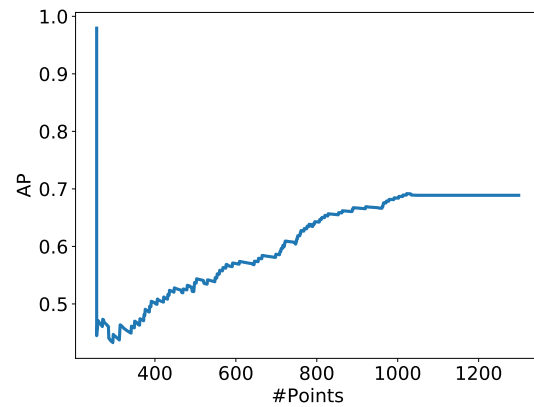
### Algorithms

- RS-Hash: We use decay constant to be 0.015, and use first 256 samples to setup the 100 components, and then compute AUC/AP over time.
- HS-Trees:
- LODA:

### Results



(a) Low/Mid - dimensional dataset analysis



(b) High dimensional dataset analysis

Figure 17: AP for two of the baseline datasets over time, using RS-Hash Streaming Algorithm.

## References

- Ahle, Thomas D et al. (2017). "Parameter-free locality sensitive hashing for spherical range reporting". In: *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*. Society for Industrial and Applied Mathematics, pp. 239–256.
- Andoni, Alexandr et al. (2017). "LSH forest: Practical algorithms made theoretical". In: *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, pp. 67–78.
- Aumasson, Jean-Philippe et al. (2012). "SipHash: A Fast Short-Input PRF." In: *INDOCRYPT*. Vol. 7668. Springer, pp. 489–508.
- Bawa, Mayank et al. (2005). "LSH forest: self-tuning indexes for similarity search". In: *Proceedings of the 14th international conference on World Wide Web*. ACM, pp. 651–660.
- Datar, Mayur et al. (2004). "Locality-sensitive hashing scheme based on p-stable distributions". In: *Proceedings of the twentieth annual symposium on Computational geometry*. ACM, pp. 253–262.
- Li, Ping et al. (2006). "Very sparse random projections". In: *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, pp. 287–296.