RPS: A Generic Reservoir Patterns Sampler

Lamine DIOP*, Marc PLANTEVIT*, Arnaud SOULET+

(*) EPITA Research Laboratory (LRE), France

(+) University of Tours, LIFAT, France

Combinatorics and Life Sciences
University of Lyon 1
October 2, 2025

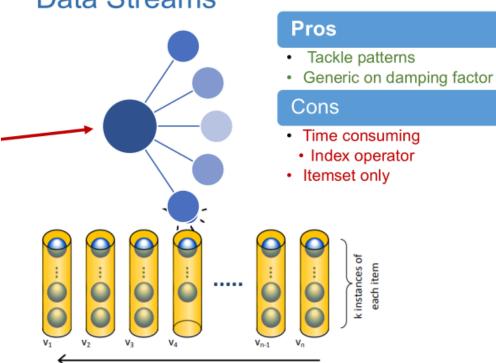
Reservoir Pattern Sampling Problem

ResPat: Reservoir Patterns Sampling in WRS-k-W: Weighted Random **Data Streams** Sampling over Data Streams Pros Tackle patterns · Generic on damping factor Cons · Time consuming · Index operator · Itemset only [Pavlos S. Efraimidis, 2015] [Giacometti & Soulet, ECMLPKDD'22]

Batch-based Reservoir Pattern Sampling

ResPat: Reservoir **Patterns** Sampling in Data Streams

[Giacometti & Soulet, ECMLPKDD'22]



RPS: Our contribution

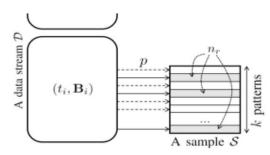


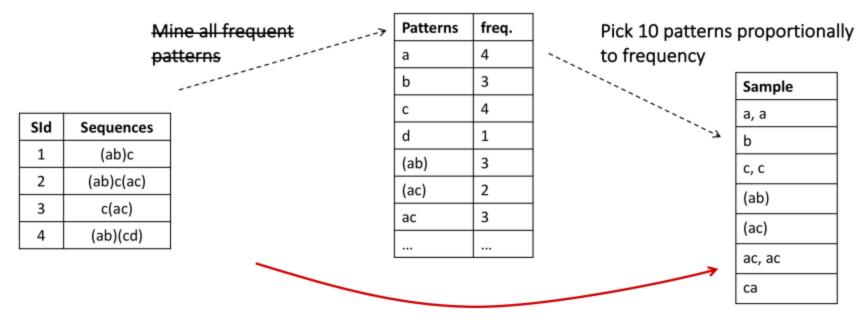
Fig. 1: Overview of the approach (the incomplete block denotes next batch)

Property 2 (From CBPD to IBF). Let k be the total number of trials, $n_r \leq k$ the minimum number of times the event must occur, and p the probability of the event occurring in a single trial. The Cumulative Binomial Probability Distribution can be computed as follows:

$$\mathbb{P}(n_r, k) \equiv \sum_{i=1}^{k} {k \choose j} p^j (1-p)^{k-j} = \mathbf{I}_p(n_r, k-n_r+1).$$
 (2)

Direct Local Pattern Sampling

Frequent sequential pattern sampling



Goal: find directly a sample

Pattern sampling state-of-the-art

Lang	juage		
Graphs	[Hasan & Zaki, 2009]		[Diop et Plantevit, IEEEBigData'24]
Sequences			[Diop et al., ICDM'2018] [Diop et al., KAIS'2019] [Diop & Ba, IEEEBigData'2021]
Transactions	[Boley <i>et al.</i> , 2010] [Li et Zaki, 2012] [Bendimerad et al.,2020]	[Dzyuba et al., 2017] [Gueguen <i>et al.</i> , 2019]	[Boley et al., SIGKDD'2011] [Boley <i>et al.</i> , 2012] [Moens & Boley, 2014] [Diop et al., ADBIS'2020] [Diop et al., IS'2022]
Numerical data			[Giacometti & Soulet, 2018] [Diop, PAKDD2022] [Djawad et al., DKE'25 (under submission
•	MCMC	SAT	Two-step random Approach

Challenges in Complex Data Structure

Challenges for subsequences

Challenge 1: compute the number of subsequences

SId	Sequences
1	(ab)c
2	(ab)c(ac)
3	c(ac)
4	(ab)(cd)

#subseq.	
6	
12	
5	
10	

10

1. Pick a sequence proportionally to its number of subsequences

Challenge 2: do not favo	0
subsequence having	
multiple occurrences	

TId	Subsequences
1	a, b, (ab), ac, bc
2	a, b, c, (ab), (ac), ac, bc,
3	c, a, (ac), ca, cc
4	a, b, c, d, (ab), (cd),

	Sample
	a, a
	b
	с, с
	(ab)
	(ac)
	ac, ac
Г	ca

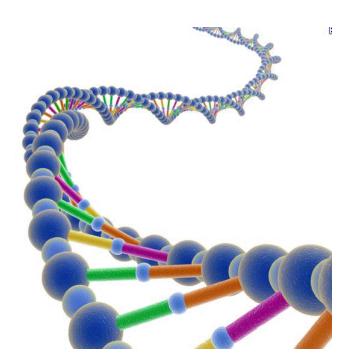
2. Pick a subsequence uniformly within the drawn sequence

Example

$$\phi_{\leq 2}(\langle (ab)c\rangle \circ (c)) \neq \phi_{\leq 2}(\langle (ab)c\rangle) \times \binom{1}{0} + \phi_{\leq 1}(\langle (ab)c\rangle) \times \binom{1}{1}$$
With repetitions, some subsequences are counted multiple times
$$\begin{pmatrix} \langle \rangle, \\ \langle a\rangle, \langle b\rangle, \langle c\rangle, \\ \langle (ab)\rangle, \langle ac\rangle, \langle bc\rangle \end{pmatrix} \circ \emptyset \qquad \begin{pmatrix} \langle \rangle, \\ \langle a\rangle, \langle b\rangle, \langle c\rangle \end{pmatrix} \circ (c)$$

$$\neq 7 \times 1 \qquad + 4 \times 1$$

$$= 11 - (-1)^2 \phi_{\leq 1}(\langle (ab)\rangle \circ (c)) = 11 - 3 = 8$$



DNA sequential data

[Diop et al., ICDM'2018]

Theorem for Distinct Subsequences Counting

Theorem 1 (Subsequence number with a maximum norm): Let s be a sequence, Y be an itemset and j be an integer, the number of distinct subsequences having a norm less or equal to j in $s \circ Y$, denoted by $\Phi_{\leq j}(s \circ Y)$, is defined as follows¹:

$$\Phi_{\leq j}(s \circ Y) = \left(\sum_{k=0}^{j} \Phi_{\leq j-k}(s) \times \binom{|Y|}{k}\right) - R_{\leq j}(s, Y)$$

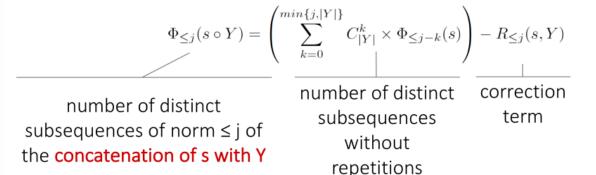
where $R_{\leq j}(s, Y)$ is the correction term defined by:

$$R_{\leq j}(s,Y) = \sum_{\emptyset \subset K \subset L(s,Y)} (-1)^{|K|+1} R_{\leq j}^K(s,Y)$$

with
$$R_{\leq j}^K(s,Y)=\sum_{k=1}^j\Phi_{\leq j-k}(s^{min(K)-1})\times {|s[K]\cap Y|\choose k}$$
 where $s[K]=\cap_{k\in K}s[k]$.

Using the inclusion-exclusion principle

1. Number of subsequences with a norm ≤ j



2. Drawing uniformly a subsequence

- □ Algorithm for drawing a subsequence within s:
 - 1. Draw randomly an occurrence o from s
 - 2. If o is not the first occurrence in s of its subsequence, repeat step 1.
 - 3. Return the subsequence stemming from o

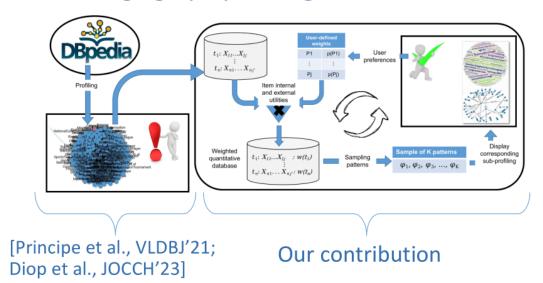
Example: 2 occurrences of **ac** in **(ab)c(ac)**:

(ab)c(ac)
$$\frac{(ab)c(ac)}{(ab)}$$

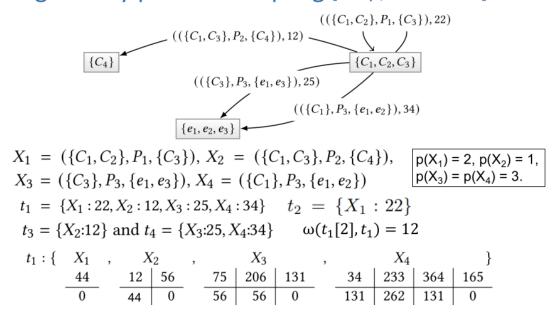
Pattern Counting in Weighted Itemsets (1/2)

High utility pattern mining & Knowledge graph profiling

Knowledge graph profiling [Principe et al, VLDBJ'21]



High utility pattern sampling [Diop, PAKDD'22]



Pattern Counting in Weighted Itemsets (2/2)

1. Weighting approach

DEFINITION 3 (UPPER TRIANGLE UTILITY). Given a quantitative transaction t, ℓ and i two positive integers less than or equal to |t|. The upper triangle utility of t denoted by V_t is defined as follows:

$$\mathcal{V}_t(\ell, i) = 0 \text{ if } \ell > i$$
 $\mathcal{V}_t(1, i) = \sum_{j=1}^i \omega(t[j], t)$

$$\mathcal{V}_t(\ell,i) = \binom{i-1}{\ell-1} \times \omega(t[i],t) + \mathcal{V}_t(\ell-1,i-1) + \mathcal{V}_t(\ell,i-1).$$

Theorem 1. Let V_t be the upper triangle utility of any given a quantitative transaction t. For any positive integers ℓ and i such that $\ell \leq i \leq |t|$, the following statement holds: $V_t(\ell, i) = \binom{i-1}{\ell-1} \times V_t(1, i)$.

PROPERTY 1 (SYMMETRY). If V_t is the upper triangle utility of a given transaction t, and ℓ and i two positive integers such that $\ell \leq i \leq |t|$ then we have: $V_t(\ell, i) = V_t(i - \ell + 1, i)$.

Theorem 2. Let t be a transaction from a qDB. The cumulative utility sum of the entire set of patterns that appear in transaction t is expressed as: $W(t) = 2^{|t|-1} \times \sum_{X \in t} \omega(X, t)$.

2. Drawing approach

Property 2 (Dichotomous random search). Let t be a quantitative transaction, and t[j] be the last picked item during the drawing process. Using a sequential random variable generation to draw the next item t[i], with i < j, to add in the sampled pattern φ is equivalent to jumping directly on it based on dichotomous search. In other words, if a random number α_i drawn from the interval $]0, \mathcal{V}_t(\ell,i) + \binom{i}{\ell} \times \mathbb{U}_{\mathcal{L}}(\varphi,t)]$ allows us to pick t[i], i.e., $\frac{\alpha_i}{\mathcal{V}_t(\ell,i) + \binom{i}{\ell} \times \mathbb{U}_{\mathcal{L}}(\varphi,t)} \leq \mathbb{P}(t[i]|\varphi \wedge \ell)$, then from position j of t, we can directly select the item t[i] such that $\mathcal{V}_t(\ell,i-1) + \binom{i-1}{\ell} \times \mathbb{U}_{\mathcal{L}}(\varphi,t) < \alpha_i \leq \mathcal{V}_t(\ell,i) + \binom{i-1}{\ell-1} \times \mathbb{U}_{\mathcal{L}}(\varphi,t)$.

Generic Problem Reformulation

Definition: (Pattern). $\varphi = \langle X'_1, \dots, X'_{n'} \rangle$ is a pattern or an generalization of an instance $\gamma = \langle X_1, \dots, X_n \rangle$, denoted by $\varphi \leq \gamma$, if there exists an index sequence $1 \leq i_1 < i_2 < \dots < i_{n'} \leq n$ such that for all $j \in [1..n']$, one has $X'_j \subseteq X_{i_j}$.

This definition is usually used in the context of sequential pattern mining, but we recall that an itemset is nothing else that a sequential pattern of length 1.

Algorithm 1 Sample_{n_r}($\mathcal{L}, \mathbf{B}, m$)

Output: A sample of n_r patterns of $\mathcal{L}(\mathbf{B})$ drawn proportionally to the utility measure m

- 1: Let $\omega_m(\gamma) \leftarrow \sum_{\varphi \prec \gamma} m(\varphi, \gamma)$ for all $\gamma \in \mathbf{B}$
- 2: Let $\phi \leftarrow \emptyset$
- 3: **for** $j \in [1..n_r]$ **do**
- 4: Draw γ from **B** with $\mathbb{P}(\gamma, \mathbf{B}) = \frac{\omega_m(\gamma)}{\sum_{\gamma_i \in \mathbf{B}} \omega_m(\gamma_i)}$
- 5: Let $\omega_m^{\gamma}(\ell) \leftarrow \sum_{\varphi \preceq \gamma \land ||\varphi|| = \ell} m(\varphi, \gamma)$ for $\ell \in [1..||\gamma||]$
- 6: Draw an integer ℓ' with a probability of $\frac{\omega_m^{\gamma}(\ell')}{\sum_{\ell} \omega_m^{\gamma}(\ell)}$
- 7: Let $\varphi_j \sim m(\{\varphi \leq \gamma : \|\varphi\| = \ell'\})$
- 8: $\phi \leftarrow \phi \cup \{\varphi_j\}$
- 9: return ϕ

RPS: Generic Reservoir Pattern Sampler

Algorithm 2 RPS: A Generic Stream pattern sampler

```
Input: A data stream \mathcal{D}, a utility measure m, a damping factor \varepsilon \in [0,1], and the desired reservoir size k
Output at time t_n: A sample \mathcal{S} of k patterns drawn in \mathcal{L}(\mathcal{D} = \langle (t_1, \mathbf{B}_1), \dots, (t_n, \mathbf{B}_n) \rangle) based on m and \varepsilon
```

- 1: $\mathcal{S} \leftarrow \emptyset$; $\mathbf{Z}_0 = 0$
- 2: **while** (t_i, \mathbf{B}_i) is from \mathcal{D} **do**

//Batch acceptance probability

```
3: \mathbf{Z}_i = \mathbf{Z}_{i-1} + \omega_m(\mathbf{B}_i) \times e^{\varepsilon \times t_i}
```

- 4: $p \leftarrow \frac{\omega_m(\mathbf{B}_i) \times e^{\varepsilon \times t_i}}{\mathbf{Z}_i}$
- 5: $x \leftarrow random(0,1)$
- 6: if p > x then

//Number of realisations

7:
$$n_r \leftarrow 1 + \arg[\mathbf{I}_p(n_r, k - n_r) = x] \triangleright \text{Definition } 10$$
//Patterns selection

- 8: $E \leftarrow getPatternsToRemove(S, n_r)$
- 9: $S \leftarrow S \setminus \{S[j] : \text{ for } j \in E\}$
- 10: **for** $\varphi_j \in \text{Sample}_{n_r}(\mathcal{L}, \mathbf{B}_i, m)$ **do**
- 11: $\mathcal{S} \leftarrow \mathcal{S} \cup \{(t_i, \varphi_j)\}$

Definition: (Inverse Incomplete Beta Function). The Inverse Incomplete Beta Function allows for the approximation of the number of successful trials n_r out of k trials that matches the CBPD with a given probability $x \in [0,1]$ as follows:

$$n_r = \arg_{n_r} [\mathbf{I}_p(n_r, k - n_r + 1) = x].$$

Computing Number of Realizations (1/2)

- Context: Selecting patterns for a reservoir of size k.
- Goal: Draw a set of n_r patterns proportionally to their weights from the current batch B_i , with $n_r \leq k$.
- Method:
 - Use n_r copies of batch B_i to simulate independent Bernoulli trials.
 - Simplify probability computation using Cumulative Binomial Probability Distribution (CBPD):

$$\mathbb{P}(n_r,k) \equiv \sum_{j=n_r}^k \binom{k}{j} p^j (1-p)^{k-j}$$

with $n_r \leq k$ the minimum number of times the event must occur, and p the probability of the event occurring in a single trial.

• Efficiency: Leverage the Incomplete Beta Function (IBF) to reduce computational overhead.

Computing Number of Realizations (2/2)

Incomplete Beta Function (IBF):

$$I_{x}(a,b) = \frac{1}{B(a,b)} \int_{0}^{x} t^{a-1} (1-t)^{b-1} dt,$$

where B(a, b) is the beta function.

- Key Application:
 - Replace summation over binomial probabilities with IBF.
 - Efficiently compute probabilities for large reservoirs.
- Inverse IBF:
 - Determine n_r for a given x:

$$n_r = \arg_{n_r} [I_p(n_r, k - n_r + 1) = x].$$

Total number of realizations

Let x = random(0, 1) and $p = \frac{\omega_m(B_i) \cdot e^{\varepsilon \cdot t_i}}{Z_i}$ for a batch current batch B_j . The total number of realizations can be computed as follows:

$$n_r = 1_{(x \leq p)} \cdot ig(1 + \operatorname{arg}_{n_r} \left[\mathsf{I}_p(n_r, k - n_r) = x
ight] ig)$$

Soundness of RPS

For any Two-Step pattern sampler, regardless the data structure

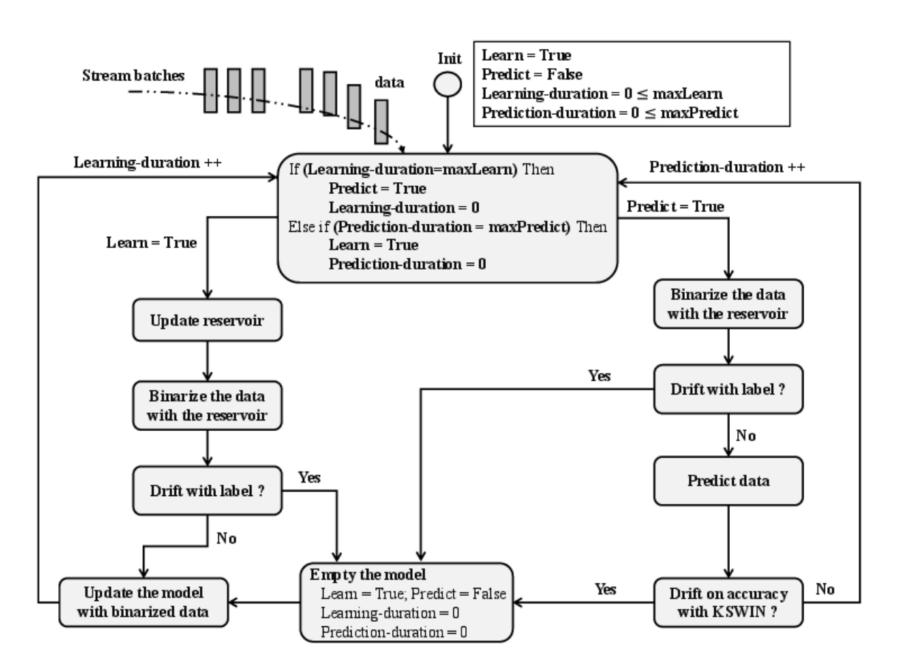
Property 5 (Soundness). Let $\mathcal{D} = \langle (t_1, \mathbf{B}_1), \dots, (t_n, \mathbf{B}_n) \rangle$ be a stream data defined over a pattern language \mathcal{L} , m be a norm-based utility measure, k the size of the reservoir \mathcal{S} and $\varepsilon \in [0,1]$ a damping factor. After observing (t_n, \mathbf{B}_n) , RPS returns a sample of k patterns $\varphi_1, \dots, \varphi_k$ where each pattern $\varphi_j = \mathcal{S}[j]$ is drawn with a probability equals to:

$$\mathbb{P}(\varphi_j = \mathcal{S}[j], \mathcal{D}) = \frac{G_m^{\varepsilon}(\varphi_j, \mathcal{D})}{\sum_{\varphi \in \mathcal{L}} G_m^{\varepsilon}(\varphi_j, \mathcal{D})}.$$

Definition: (Pattern Global Utility under temporal bias). Let $\mathcal{D} = \langle (t_1, \mathbf{B}_1), \dots, (t_n, \mathbf{B}_n) \rangle$ be a stream data defined over a pattern language \mathcal{L} , m be an interestingness utility measure and $\varepsilon \in [0,1]$ a damping factor. At time t_n , the global utility of any pattern φ inserted into the reservoir at time t and that still appears into \mathcal{S} is given by:

$$G_m^{\varepsilon}(\varphi, \mathcal{D}) = \sum_{(t_i, \mathbf{B}_i) \in \mathcal{D}} \left(\left(\sum_{\gamma_{i_j} \in \mathbf{B}_i} m(\varphi, \gamma_{i_j}) \right) \times \nabla_{\varepsilon}(t_n, t_i) \right).$$

Use Case on Classification of Stream
Sequential Data



Snapshots of RPS-Classifiers Evolution per batch

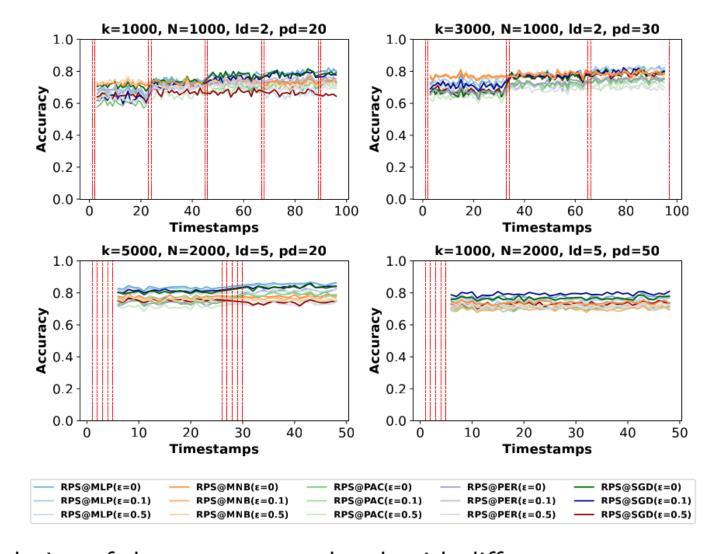


Figure: Evolution of the accuracy per batch with different parameters on Books. The learning timestamps are colored red. k: reservoir size, N: batch size, Id: learning duration, pd: predict duration

Online Classifier Accuracies

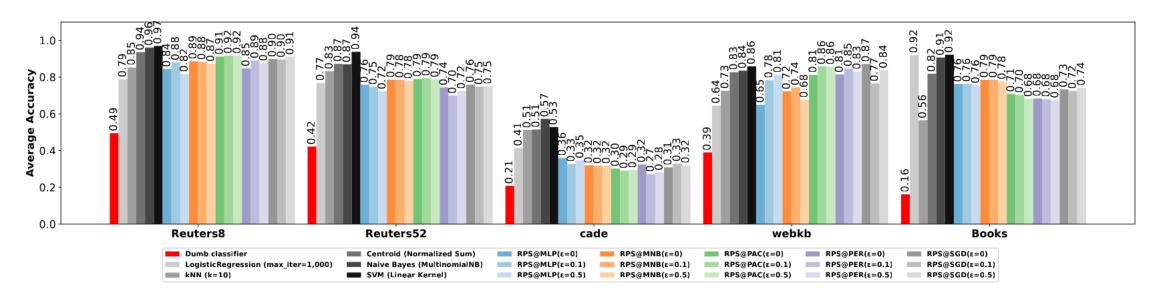


Figure: Comparison between RPS-based classifiers (with reservoir size k = 10,000; batch size=1,000; learning duration=2 time-units, predict duration=52 time-units) vs cheater classifiers (with 50% train and 50% test)

Conclusion

Summary:

- Introduced a novel weighted reservoir sampling approach.
- Demonstrated scalability and effectiveness in complex streams.

Impact:

- Enhances scalability for large-scale systems.
- Reduces computational complexity for real-time tasks.

Future Work:

- Extension to multi-stream environments.
- For xAI: Capturing patterns from latent space during model training