

# Verification of the CVM algorithm with a New Recursive Analysis Technique

Emin Karayel, Derek Khu, Kuldeep S. Meel, Yong Kiam Tan,  
and Seng Joe Watt

March 19, 2025

## Abstract

In 2022, Chakraborty et al. [1] published a streaming algorithm (henceforth, the CVM algorithm) for the distinct elements problem, that deviated considerably from the state-of-the-art, due to its simplicity and avoidance of standard derandomization techniques, while still maintaining a close to optimal logarithmic space complexity.

In this entry, we verify the CVM algorithm's correctness using a new technique which simplifies the analysis considerably compared to the original proof by Chakraborty et al. The main idea is based on a probabilistic invariant that allows us to derive concentration bounds using the Cramér–Chernoff method.

This new technique opens up the possible algorithm design space, and we introduce a new variant of the CVM algorithm, that is total, and also has an additional property in addition to concentration: unbiasedness. This means the expected result of the algorithm is exactly equal to the desired result. The latter is also a new property, that neither the original CVM algorithm nor classic algorithms for the distinct elements problem possess.

# Contents

<b>1</b>	<b>Preliminary Definitions and Results</b>	<b>3</b>
<b>2</b>	<b>Abstract Algorithm</b>	<b>5</b>
<b>3</b>	<b>The Original CVM Algorithm</b>	<b>12</b>
<b>4</b>	<b>The New Unbiased Algorithm</b>	<b>15</b>
<b>A</b>	<b>Informal Proof</b>	<b>19</b>
	A.1 Loop Invariant . . . . .	19
	A.2 Concentration . . . . .	22
	A.3 Unbiasedness . . . . .	24

# 1 Preliminary Definitions and Results

```

theory CVM-Preliminary
  imports HOL-Probability.SPMF
begin

lemma bounded-finite:
  assumes ⟨finite S⟩
  shows ⟨bounded (f ' S)⟩
  ⟨proof⟩

lemma of-bool-power:
  assumes ⟨y > 0⟩
  shows ⟨(of-bool x::real) ^ (y::nat) = of-bool x⟩
  ⟨proof⟩

lemma card-filter-mono:
  assumes ⟨finite S⟩
  shows ⟨card (Set.filter p S) ≤ card S⟩
  ⟨proof⟩

fun foldM ::
  ⟨('a ⇒ ('b ⇒ 'c) ⇒ 'c) ⇒ ('b ⇒ 'c) ⇒ ('d ⇒ 'b ⇒ 'a) ⇒ 'd list ⇒ 'b ⇒ 'c) where
  ⟨foldM - return' - [] val = return' val⟩ |
  ⟨foldM bind' return' f (x # xs) val =
    bind' (f x val) (foldM bind' return' f xs)⟩

abbreviation foldM-pmf ::
  ⟨('a ⇒ 'b ⇒ 'b pmf) ⇒ 'a list ⇒ 'b ⇒ 'b pmf⟩ where
  ⟨foldM-pmf ≡ foldM bind-pmf return-pmf⟩

lemma foldM-pmf-snoc: ⟨foldM-pmf f (xs@[y]) val = bind-pmf (foldM-pmf f xs val) (f y)⟩
  ⟨proof⟩

abbreviation foldM-spmf
  :: ⟨('a ⇒ 'b ⇒ 'b spmf) ⇒ 'a list ⇒ 'b ⇒ 'b spmf⟩ where
  ⟨foldM-spmf ≡ foldM bind-spmf return-spmf⟩

lemma foldM-spmf-snoc: ⟨foldM-spmf f (xs@[y]) val = bind-spmf (foldM-spmf f xs val) (f y)⟩
  ⟨proof⟩

abbreviation ⟨prob-fail ≡ (λx. pmf x None)⟩

abbreviation ⟨fail-spmf ≡ return-pmf None⟩

abbreviation fails-or-satisfies :: ⟨('a ⇒ bool) ⇒ 'a option ⇒ bool⟩ where
  ⟨fails-or-satisfies ≡ case-option True⟩

lemma prob-fail-foldM-spmf-le :
  fixes
  p :: real and
  P :: ⟨'b ⇒ bool⟩ and

```

$f :: \langle 'a \Rightarrow 'b \Rightarrow 'b \text{ spmf} \rangle$   
**assumes**  
 $\langle \bigwedge x y z. P y \Longrightarrow z \in \text{set-spmf } (f x y) \Longrightarrow P z \rangle$   
 $\langle \bigwedge x \text{ val}. P \text{ val} \Longrightarrow \text{prob-fail } (f x \text{ val}) \leq p \rangle$   
 $\langle P \text{ val} \rangle$   
**shows**  $\langle \text{prob-fail } (\text{foldM-spmf } f \text{ xs val}) \leq \text{real } (\text{length } \text{xs}) * p \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *foldM-spmf-of-pmf-eq* :  
**shows**  $\langle \text{foldM-spmf } (\lambda x y. \text{spmf-of-pmf } (f x y)) \text{ xs} = \text{spmf-of-pmf} \circ \text{foldM-pmf } f \text{ xs} \rangle$   
**(is ?thesis-0)**  
**and**  $\langle \text{foldM-spmf } (\lambda x y. \text{spmf-of-pmf } (f x y)) \text{ xs val} = \text{spmf-of-pmf } (\text{foldM-pmf } f \text{ xs val}) \rangle$   
**(is ?thesis-1)**  
 $\langle \text{proof} \rangle$

**end**

## 2 Abstract Algorithm

This section verifies an abstract version of the CVM algorithm, where the subsampling step can be an arbitrary randomized algorithm fulfilling an expectation invariant.

The abstract algorithm is presented in Algorithm 1.

---

**Algorithm 1** Abstract CVM algorithm.

---

**Input:** Stream elements  $a_1, \dots, a_l$ ,  $0 < \varepsilon$ ,  $0 < \delta < 1$ ,  $\frac{1/2}{\varepsilon} f < 1$

**Output:** An estimate  $R$ , s.t.,  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$  where  $A := \{a_1, \dots, a_l\}$ .

```

1:  $\chi \leftarrow \{\}, p \leftarrow 1, n \geq \lceil \frac{12}{\varepsilon^2} \ln(\frac{3l}{\delta}) \rceil$ 
2: for  $i \leftarrow 1$  to  $l$  do
3:    $b \stackrel{\$}{\leftarrow} \text{Ber}(p)$  ▷ insert  $a_i$  with probability  $p$  (and remove it otherwise)
4:   if  $b$  then
5:      $\chi \leftarrow \chi \cup \{a_i\}$ 
6:   else
7:      $\chi \leftarrow \chi - \{a_i\}$ 
8:   if  $|\chi| = n$  then
9:      $\chi \stackrel{\$}{\leftarrow} \text{subsample}(\chi)$  ▷ abstract subsampling step
10:     $p \leftarrow pf$ 
11: return  $\frac{|\chi|}{p}$  ▷ estimate cardinality of  $A$ 

```

---

For the subsampling step we assume that it fulfills the following inequality:

$$\int_{\text{subsample}(\chi)} \left( \prod_{i \in S} g(i \in \omega) \right) d\omega \leq \prod_{i \in S} \left( \int_{\text{Ber}(f)} g(\omega) d\omega \right) \quad (1)$$

for all non-negative functions  $g$  and  $S \subseteq \chi$ , where  $\text{Ber}(p)$  denotes the Bernoulli-distribution.

The original CVM algorithm uses a subsampling step where each element of  $\chi$  is retained independently with probability  $f$ . It is straightforward to see that this fulfills the above condition (with equality).

The new CVM algorithm variant proposed in this work uses a subsampling step where a random  $nf$ -sized subset of  $\chi$  is kept. This also fulfills the above inequality, although this is harder to prove and will be explained in more detail in Section 4.

In this section, we will verify that the above abstract algorithm indeed fulfills the desired conditions on its estimate, as well as unbiasedness, i.e., that:  $\mathbb{E}[R] = |A|$ . The part that is not going to be verified in this section, is the fact that the algorithm keeps at most  $n$  elements in the state  $\chi$ , because it is not unconditionally true, but will be ensured (by different means) for the concrete instantiations in the following sections.

An informal version of this proof is presented in Appendix A. For important lemmas and theorems, we include a reference to the corresponding statement in the appendix.

**theory** *CVM-Abstract-Algorithm*

```

imports
  HOL-Decision-Procs.Approximation
  CVM-Preliminary
  Finite-Fields.Finite-Fields-More-PMF
  Universal-Hash-Families.Universal-Hash-Families-More-Product-PMF
begin

unbundle no vec-syntax

datatype 'a state = State (state- $\chi$ : 'a set) (state-p: real)

datatype 'a run-state = FinalState 'a list | IntermState 'a list 'a

lemma run-state-induct:
  assumes 'P (FinalState [])
  assumes 'P ( $\bigwedge xs x. P$  (FinalState xs)  $\implies P$  (IntermState xs x))
  assumes 'P ( $\bigwedge xs x. P$  (IntermState xs x)  $\implies P$  (FinalState (xs@[x])))
  shows 'P result
  (proof)

locale cvm-algo-abstract =
  fixes n :: nat and f :: real and subsample :: 'a set  $\Rightarrow$  'a set pmf
  assumes n-gt-0: <n > 0>
  assumes f-range: <f  $\in$  {1/2.. $<1$ }>
  assumes subsample:
    < $\bigwedge U x. \text{card } U = n \implies x \in \text{set-pmf } (\text{subsample } U) \implies x \subseteq U$ >
  assumes subsample-inequality:
    < $\bigwedge g U S. S \subseteq U$ 
       $\implies \text{card } U = n$ 
       $\implies \text{range } g \subseteq \{0::\text{real}..\}$ 
       $\implies (\int \omega. (\prod_{s \in S}. g(s \in \omega)) \partial \text{subsample } U) \leq (\prod_{s \in S}. (\int \omega. g \omega \partial \text{bernoulli-pmf } f))$ >
begin

```

Line 1 of Algorithm 1:

```

definition initial-state :: 'a state where
  <initial-state  $\equiv$  State {} 1>

```

Lines 3–7:

```

fun step-1 :: 'a  $\Rightarrow$  'a state  $\Rightarrow$  'a state pmf where
  <step-1 a (State  $\chi$  p) =
    do {
      b  $\leftarrow$  bernoulli-pmf p;
      let  $\chi = (\text{if } b \text{ then } \chi \cup \{a\} \text{ else } \chi - \{a\});$ 

      return-pmf (State  $\chi$  p)
    }>

```

Lines 8–10:

```

fun step-2 :: 'a state  $\Rightarrow$  'a state pmf where
  <step-2 (State  $\chi$  p) = do {
    if card  $\chi = n$ 
    then do {

```

```

     $\chi \leftarrow \text{subsample } \chi;$ 
     $\text{return-pmf } (\text{State } \chi (p*f))$ 
  } else do {
     $\text{return-pmf } (\text{State } \chi p)$ 
  }
}
```

**schematic-goal** *step-1-def*:  $\langle \text{step-1 } x \sigma = ?x \rangle$   
 $\langle \text{proof} \rangle$

**schematic-goal** *step-2-def*:  $\langle \text{step-2 } \sigma = ?x \rangle$   
 $\langle \text{proof} \rangle$

Lines 1–10:

**definition** *run-steps* ::  $\langle 'a \text{ list} \Rightarrow 'a \text{ state pmf} \rangle$  **where**  
 $\langle \text{run-steps } xs \equiv \text{foldM-pmf } (\lambda x \sigma. \text{step-1 } x \sigma \gg \text{step-2}) \text{ xs initial-state} \rangle$

Line 11:

**definition** *estimate* ::  $\langle 'a \text{ state} \Rightarrow \text{real} \rangle$  **where**  
 $\langle \text{estimate } \sigma = \text{card } (\text{state-}\chi \sigma) / \text{state-p } \sigma \rangle$

**lemma** *run-steps-snoc*:  $\langle \text{run-steps } (xs @ [x]) = \text{run-steps } xs \gg \text{step-1 } x \gg \text{step-2} \rangle$   
 $\langle \text{proof} \rangle$

**fun** *run-state-pmf* **where**  
 $\langle \text{run-state-pmf } (\text{FinalState } xs) = \text{run-steps } xs \rangle$  |  
 $\langle \text{run-state-pmf } (\text{IntermState } xs x) = \text{run-steps } xs \gg \text{step-1 } x \rangle$

**fun** *len-run-state* **where**  
 $\langle \text{len-run-state } (\text{FinalState } xs) = \text{length } xs \rangle$  |  
 $\langle \text{len-run-state } (\text{IntermState } xs x) = \text{length } xs \rangle$

**fun** *run-state-set* **where**  
 $\langle \text{run-state-set } (\text{FinalState } xs) = \text{set } xs \rangle$  |  
 $\langle \text{run-state-set } (\text{IntermState } xs x) = \text{set } xs \cup \{x\} \rangle$

**lemma** *finite-run-state-set[simp]*:  $\langle \text{finite } (\text{run-state-set } \sigma) \rangle$   $\langle \text{proof} \rangle$

**lemma** *subsample-finite-pmf*:  
**assumes**  $\langle \text{card } U = n \rangle$   
**shows**  $\langle \text{finite } (\text{set-pmf } (\text{subsample } U)) \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *finite-run-state-pmf*:  $\langle \text{finite } (\text{set-pmf } (\text{run-state-pmf } \rho)) \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *state- $\chi$ -run-state-pmf*:  $\langle \text{AE } \sigma \text{ in run-state-pmf } \rho. \text{state-}\chi \sigma \subseteq \text{run-state-set } \rho \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *state- $\chi$ -finite*:  $\langle \text{AE } \sigma \text{ in run-state-pmf } \rho. \text{finite } (\text{state-}\chi \sigma) \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *state-p-range*:  $\langle \text{AE } \sigma \text{ in run-state-pmf } \rho. \text{state-p } \sigma \in \{0 <.. 1\} \rangle$

*<proof>*

Lemma 1:

**lemma** *run-steps-preserves-expectation-le*:

**fixes**  $\varphi :: \langle \text{real} \Rightarrow \text{bool} \Rightarrow \text{real} \rangle$

**assumes** *phi* :

$\langle \bigwedge x b. \llbracket 0 < x; x \leq 1 \rrbracket \Longrightarrow \varphi x b \geq 0 \rangle$

$\langle \bigwedge p x. \llbracket 0 < p; p \leq 1; 0 < x; x \leq 1 \rrbracket \Longrightarrow (\int \omega. \varphi x \omega \partial \text{bernoulli-pmf } p) \leq \varphi (x / p)$

*True* $\rangle$

$\langle \text{mono-on } \{0..1\} (\lambda x. \varphi x \text{ False}) \rangle$

**defines**  $\langle \text{aux} \equiv \lambda S \sigma. (\prod x \in S. \varphi (\text{state-p } \sigma) (x \in \text{state-}\chi \sigma)) \rangle$

**assumes**  $\langle S \subseteq \text{run-state-set } \varrho \rangle$

**shows**  $\langle \text{measure-pmf.expectation } (\text{run-state-pmf } \varrho) (\text{aux } S) \leq \varphi 1 \text{ True} \wedge \text{card } S \rangle$

*<proof>*

Lemma 2:

**lemma** *run-steps-preserves-expectation-le'* :

**fixes**  $q :: \text{real}$  **and**  $h :: \langle \text{real} \Rightarrow \text{real} \rangle$

**assumes** *h*:

$\langle 0 < q \rangle \langle q \leq 1 \rangle$

$\langle \text{concave-on } \{0 .. 1 / q\} h \rangle$

$\langle \bigwedge x. \llbracket 0 \leq x; x * q \leq 1 \rrbracket \Longrightarrow h x \geq 0 \rangle$

**defines**

$\langle \text{aux} \equiv \lambda S \sigma. (\prod x \in S. \text{of-bool } (\text{state-p } \sigma \geq q) * h (\text{of-bool } (x \in \text{state-}\chi \sigma) / \text{state-p } \sigma)) \rangle$

**assumes**  $\langle S \subseteq \text{run-state-set } \varrho \rangle$

**shows**  $\langle (\int \tau. \text{aux } S \tau \partial \text{run-state-pmf } \varrho) \leq (h 1) \wedge \text{card } S \rangle$  (**is**  $\langle ?L \leq ?R \rangle$ )

*<proof>*

Analysis of the case where  $n \leq \text{card } (\text{set } xs)$ :

**context**

**fixes**  $xs :: \langle 'a \text{ list} \rangle$

**begin**

**private abbreviation**  $\langle A \equiv \text{real } (\text{card } (\text{set } xs)) \rangle$

**context**

**assumes** *set-larger-than-n*:  $\langle \text{card } (\text{set } xs) \geq n \rangle$

**begin**

**private definition**  $\langle q = \text{real } n / (4 * \text{card } (\text{set } xs)) \rangle$

**lemma** *q-range*:  $\langle q \in \{0 < .. 1/4\} \rangle$

*<proof>*

**lemma** *mono-nonnegI*:

**assumes**  $\langle \bigwedge x. x \in I \Longrightarrow h' x \geq 0 \rangle$

**assumes**  $\langle \bigwedge x. x \in I \Longrightarrow (h \text{ has-real-derivative } (h' x)) (\text{at } x) \rangle$

**assumes**  $\langle x \in I \cap \{0..\} \rangle \langle \text{convex } I \rangle \langle 0 \in I \rangle \langle h 0 \geq 0 \rangle$

**shows**  $\langle h x \geq 0 \rangle$

*<proof>*



**lemma** *upper-tail-bound-helper*:

**assumes**  $\langle x \in \{0 <.. 1 :: \text{real}\} \rangle$

**defines**  $\langle h \equiv (\lambda x. -q * x^2 / 3 - \ln(1 + q * x) + q * \ln(1 + x) * (1 + x)) \rangle$

**shows**  $\langle h x \geq 0 \rangle$

$\langle \text{proof} \rangle$  **definition**  $\vartheta$  **where**  $\langle \vartheta t x = 1 + q * x * (\exp(t / q) - 1) \rangle$

**lemma**  *$\vartheta$ -concave*:  $\langle \text{concave-on } \{0..1 / q\} (\vartheta t) \rangle$

$\langle \text{proof} \rangle$

**lemma**  *$\vartheta$ -ge-exp-1*:

**assumes**  $\langle x \in \{0..1/q\} \rangle$

**shows**  $\langle \exp(t * x) \leq \vartheta t x \rangle$

$\langle \text{proof} \rangle$

**lemma**  *$\vartheta$ -ge-exp*:

**assumes**  $\langle y \geq q \rangle$

**shows**  $\langle \exp(t / y) \leq \vartheta t (1 / y) \rangle$

$\langle \text{proof} \rangle$

**lemma**  *$\vartheta$ -nonneg*:

**assumes**  $\langle x \in \{0..1/q\} \rangle$

**shows**  $\langle \vartheta t x \geq 0 \rangle \langle \vartheta t x > 0 \rangle$

$\langle \text{proof} \rangle$

**lemma**  *$\vartheta$ -0*:  $\langle \vartheta t 0 = 1 \rangle \langle \text{proof} \rangle$

**lemma** *tail-bound-aux*:

**assumes**  $\langle \text{run-state-set } \varrho \subseteq \text{set } xs \rangle \langle c > 0 \rangle$

**defines**  $\langle A' \equiv \text{real}(\text{card}(\text{run-state-set } \varrho)) \rangle$

**shows**  $\langle \text{measure}(\text{run-state-pmf } \varrho) \{ \omega. \exp(t * \text{estimate } \omega) \geq c \wedge \text{state-p } \omega \geq q \} \leq \vartheta t$   
 $1 \text{ pour } A'/c \rangle$

$(\text{is } \langle ?L \leq ?R \rangle)$

$\langle \text{proof} \rangle$

**Lemma 3:**

**lemma** *upper-tail-bound*:

**assumes**  $\langle \varepsilon \in \{0 <.. 1 :: \text{real}\} \rangle$

**assumes**  $\langle \text{run-state-set } \varrho \subseteq \text{set } xs \rangle$

**shows**  $\langle \text{measure}(\text{run-state-pmf } \varrho) \{ \omega. \text{estimate } \omega \geq (1 + \varepsilon) * A \wedge \text{state-p } \omega \geq q \} \leq$   
 $\exp(-\text{real } n / 12 * \varepsilon^2) \rangle$

$(\text{is } \langle ?L \leq ?R \rangle)$

$\langle \text{proof} \rangle$

**Lemma 4:**

**lemma** *low-p*:

**shows**  $\langle \text{measure}(\text{run-steps } xs) \{ \sigma. \text{state-p } \sigma < q \} \leq \text{real}(\text{length } xs) * \exp(-\text{real } n / 12) \rangle$

$(\text{is } \langle ?L \leq ?R \rangle)$

$\langle \text{proof} \rangle$

**lemma** *lower-tail-bound-helper*:

**assumes**  $\langle x \in \{0 <.. < 1 :: \text{real}\} \rangle$

**defines**  $\langle h \equiv (\lambda x. -q * x^2 / 2 - \ln(1 - q * x) + q * \ln(1 - x) * (1 - x)) \rangle$

**shows**  $\langle h x \geq 0 \rangle$

*<proof>*

Lemma 5:

**lemma** *lower-tail-bound*:

**assumes**  $\langle \varepsilon \in \{0 < .. < 1 :: \text{real}\} \rangle$

**shows**  $\langle \text{measure (run-steps } xs) \{ \omega. \text{ estimate } \omega \leq (1 - \varepsilon) * A \wedge \text{state-p } \omega \geq q \} \leq \exp(-\text{real } n / 8 * \varepsilon^2) \rangle$

(**is**  $\langle ?L \leq ?R \rangle$ )

*<proof>*

**lemma** *correctness-aux*:

**assumes**  $\langle \varepsilon \in \{0 < .. < 1 :: \text{real}\} \rangle \langle \delta \in \{0 < .. < 1 :: \text{real}\} \rangle$

**assumes**  $\langle \text{real } n \geq 12 / \varepsilon^2 * \ln (3 * \text{real (length } xs) / \delta) \rangle$

**shows**  $\langle \text{measure (run-steps } xs) \{ \omega. |\text{estimate } \omega - A| > \varepsilon * A \} \leq \delta \rangle$

(**is**  $\langle ?L \leq ?R \rangle$ )

*<proof>*

**end**

**lemma** *deterministic-phase*:

**assumes**  $\langle \text{card (run-state-set } \sigma) < n \rangle$

**shows**  $\langle \text{run-state-pmf } \sigma = \text{return-pmf (State (run-state-set } \sigma) 1) \rangle$

*<proof>*

Theorem 1:

**theorem** *correctness*:

**fixes**  $\varepsilon \delta :: \text{real}$

**assumes**  $\langle \varepsilon \in \{0 < .. < 1\} \rangle \langle \delta \in \{0 < .. < 1\} \rangle$

**assumes**  $\langle \text{real } n \geq 12 / \varepsilon^2 * \ln (3 * \text{real (length } xs) / \delta) \rangle$

**shows**  $\langle \text{measure (run-steps } xs) \{ \omega. |\text{estimate } \omega - A| > \varepsilon * A \} \leq \delta \rangle$

*<proof>*

**lemma** *p-pos*:  $\langle \exists M \in \{0 < .. 1\}. \text{AE } \omega \text{ in run-steps } xs. \text{state-p } \omega \geq M \rangle$

*<proof>*

**lemma** *run-steps-expectation-sing*:

**assumes**  $i: \langle i \in \text{set } xs \rangle$

**shows**  $\langle \text{measure-pmf.expectation (run-steps } xs) (\lambda \omega. \text{of-bool (} i \in \text{state-}\chi \omega) / \text{state-p } \omega) = 1 \rangle$

(**is**  $\langle ?L = - \rangle$ )

*<proof>*

Subsection A.3:

**corollary** *unbiasedness*:

**fixes**  $\sigma :: \langle 'a \text{ run-state} \rangle$

**shows**  $\langle \text{measure-pmf.expectation (run-steps } xs) \text{estimate} = \text{real (card (set } xs)) \rangle$

(**is**  $\langle ?L = - \rangle$ )

*<proof>*

**end**

**end**

end

### 3 The Original CVM Algorithm

In this section, we verify the algorithm as presented by Chakraborty et al. [1] (repeated, here, in Algorithm 2), with the following caveat:

In the original algorithm the elements are removed with probability  $f := \frac{1}{2}$  in the subsampling step. The version verified here allows for any  $f \in [\frac{1}{2}, e^{-1/12}]$ .

---

**Algorithm 2** Original CVM algorithm.

---

**Input:** Stream elements  $a_1, \dots, a_l$ ,  $0 < \varepsilon$ ,  $0 < \delta < 1$ ,  $f$  subsampling param.

**Output:** An estimate  $R$ , s.t.,  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$  where  $A := \{a_1, \dots, a_l\}$ .

```

1:  $\chi \leftarrow \{\}, p \leftarrow 1, n \geq \lceil \frac{12}{\varepsilon^2} \ln(\frac{6l}{\delta}) \rceil$ 
2: for  $i \leftarrow 1$  to  $l$  do
3:    $b \stackrel{\$}{\leftarrow} \text{Ber}(p)$   $\triangleright$  insert  $a_i$  with probability  $p$  (and remove it otherwise)
4:   if  $b$  then
5:      $\chi \leftarrow \chi \cup \{a_i\}$ 
6:   else
7:      $\chi \leftarrow \chi - \{a_i\}$ 
8:   if  $|\chi| = n$  then
9:      $\chi \stackrel{\$}{\leftarrow} \text{subsample}(\chi)$   $\triangleright$  keep each element of  $\chi$  indep. with prob.  $f$ 
10:     $p \leftarrow pf$ 
11:   if  $|\chi| = n$  then
12:     return  $\perp$ 
13: return  $\frac{|\chi|}{p}$   $\triangleright$  estimate cardinality of  $A$ 

```

---

The first step of the proof is identical to the original proof [1], where the above algorithm is approximated by a second algorithm, where lines 11–12 are removed, i.e., the two algorithms behave identically, unless the very improbable event—where the subsampling step fails to remove any elements—occurs. It is possible to show that the total variational distance between the two algorithms is at most  $\frac{\delta}{2}$ .

In the second step, we verify that the probability that the second algorithm returns an estimate outside of the desired interval is also at most  $\frac{\delta}{2}$ . This, of course, works by noticing that it is an instance of the abstract algorithm we introduced in Section 2. In combination, we conclude a failure probability of  $\delta$  for the unmodified version of the algorithm.

On the other hand, the fact that the number of elements in the buffer is at most  $n$  can be seen directly from Algorithm 2.

**theory** *CVM-Original-Algorithm*

**imports** *CVM-Abstract-Algorithm*

**begin**

**context**

**fixes**  $f :: \text{real}$

**fixes**  $n :: \text{nat}$

**assumes**  $f\text{-range}$ :  $\langle f \in \{1/2..exp(-1/12)\} \rangle$

**assumes**  $n\text{-gt-0}$ :  $\langle n > 0 \rangle$

**begin**

Line 1:

```
definition initial-state :: ⟨'a state⟩ where  
  ⟨initial-state = State {} 1⟩
```

Lines 3–7:

```
fun step-1 :: ⟨'a ⇒ 'a state ⇒ 'a state spmf⟩ where  
  ⟨step-1 a (State χ p) =  
    do {  
      b ← bernoulli-pmf p;  
      let χ = (if b then χ ∪ {a} else χ - {a});  
  
      return-spmf (State χ p)  
    }⟩
```

```
definition subsample :: ⟨'a set ⇒ 'a set spmf⟩ where  
  ⟨subsample χ =  
    do {  
      keep-in-χ ← prod-pmf χ (λ-. bernoulli-pmf f);  
      return-spmf (Set.filter keep-in-χ χ)  
    }⟩
```

Lines 8–10:

```
fun step-2 :: ⟨'a state ⇒ 'a state spmf⟩ where  
  ⟨step-2 (State χ p) =  
    do {  
      if card χ = n then do {  
        χ ← subsample χ;  
        return-spmf (State χ (p * f))  
      } else  
        return-spmf (State χ p)  
    }⟩
```

Lines 11–12:

```
fun step-3 :: ⟨'a state ⇒ 'a state spmf⟩ where  
  ⟨step-3 (State χ p) =  
    do {  
      if card χ = n  
      then fail-spmf  
      else return-spmf (State χ p)  
    }⟩
```

Lines 1–12:

```
definition run-steps :: ⟨'a list ⇒ 'a state spmf⟩ where  
  ⟨run-steps xs ≡ foldM-spmf (λx σ. step-1 x σ ≧≧ step-2 ≧≧ step-3) xs initial-state⟩
```

Line 13:

```
definition estimate :: ⟨'a state ⇒ real⟩ where  
  ⟨estimate σ = card (state-χ σ) / state-p σ⟩
```

```
definition run-algo :: ⟨'a list ⇒ real spmf⟩ where  
  ⟨run-algo xs = map-spmf estimate (run-steps xs)⟩
```

**schematic-goal** *step-1-m-def*:  $\langle \text{step-1 } x \sigma = ?x \rangle$   
 $\langle \text{proof} \rangle$

**schematic-goal** *step-2-m-def*:  $\langle \text{step-2 } \sigma = ?x \rangle$   
 $\langle \text{proof} \rangle$

**schematic-goal** *step-3-m-def*:  $\langle \text{step-3 } \sigma = ?x \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *ord-spmf-remove-step3*:  
 $\langle \text{ord-spmf } (=) (\text{step-1 } x \sigma \ggg \text{step-2} \ggg \text{step-3}) (\text{step-1 } x \sigma \ggg \text{step-2}) \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *ord-spmf-run-steps*:  
 $\langle \text{ord-spmf } (=) (\text{run-steps } xs) (\text{foldM-spmf } (\lambda x \sigma. \text{step-1 } x \sigma \ggg \text{step-2}) xs \text{ initial-state}) \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *f-range-simple*:  $\langle f \in \{1/2..<1\} \rangle$   
 $\langle \text{proof} \rangle$

Main result:

**theorem** *correctness*:  
**fixes**  $xs :: \langle 'a \text{ list} \rangle$   
**assumes**  $\langle \varepsilon \in \{0 < .. < 1\} \rangle \langle \delta \in \{0 < .. < 1\} \rangle$   
**assumes**  $\langle \text{real } n \geq 12 / \varepsilon^2 * \ln (6 * \text{real } (\text{length } xs) / \delta) \rangle$   
**defines**  $\langle A \equiv \text{real } (\text{card } (\text{set } xs)) \rangle$   
**shows**  $\langle \mathcal{P}(\omega \text{ in run-algo } xs. \text{fails-or-satisfies } (\lambda R. |R - A| > \varepsilon * A) \omega) \leq \delta \rangle$   
 $\langle \text{is } \langle ?L \leq ?R \rangle \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *space-usage*:  
 $\langle AE \sigma \text{ in measure-spmf } (\text{run-steps } xs). \text{card } (\text{state-}\chi \sigma) < n \wedge \text{finite } (\text{state-}\chi \sigma) \rangle$   
 $\langle \text{proof} \rangle$

**end**

**end**

## 4 The New Unbiased Algorithm

In this section, we introduce the new algorithm variant promised in the abstract. The main change is to replace the subsampling step of the original algorithm, which removes each element of the buffer independently with probability  $f$ . Instead, we choose a random  $nf$ -subset of the buffer (see Algorithm 3). (This means  $f, n$  must be chosen, such that  $nf$  is an integer.)

---

**Algorithm 3** New CVM algorithm.

---

**Input:** Stream elements  $a_1, \dots, a_l$ ,  $0 < \varepsilon$ ,  $0 < \delta < 1$ ,  $f$  subsampling param.

**Output:** An estimate  $R$ , s.t.,  $\mathcal{P}(|R - |A|| > \varepsilon|A|) \leq \delta$  where  $A := \{a_1, \dots, a_l\}$ .

```

1:  $\chi \leftarrow \{\}, p \leftarrow 1, n \geq \lceil \frac{12}{\varepsilon^2} \ln(\frac{3l}{\delta}) \rceil$ 
2: for  $i \leftarrow 1$  to  $l$  do
3:    $b \stackrel{\$}{\leftarrow} \text{Ber}(p)$   $\triangleright$  insert  $a_i$  with probability  $p$  (and remove it otherwise)
4:   if  $b$  then
5:      $\chi \leftarrow \chi \cup \{a_i\}$ 
6:   else
7:      $\chi \leftarrow \chi - \{a_i\}$ 
8:   if  $|\chi| = n$  then
9:      $\chi \stackrel{\$}{\leftarrow} \text{subsample}(\chi)$   $\triangleright$  Choose a random  $nf$ -subset of  $\chi$ 
10:     $p \leftarrow pf$ 
11: return  $\frac{|\chi|}{p}$   $\triangleright$  estimate cardinality of  $A$ 

```

---

The fact that this still preserves the required inequality for the subsampling operation (Eq. 1) follows from the negative associativity of permutation distributions [2, Th. 10].

(See also our formalization of the concept [3].)

Because the subsampling step always removes elements unconditionally, the second check, whether the subsampling succeeded of the original algorithm is not necessary anymore.

This improves the space usage of the algorithm, because the first reduction argument from Section 3 is now obsolete. Moreover the resulting algorithm is now unbiased, because it is an instance of the abstract algorithm of Section 2.

**theory** *CVM-New-Unbiased-Algorithm*

**imports**

*CVM-Abstract-Algorithm*

*Probabilistic-Prime-Tests.Generalized-Primality-Test*

*Negative-Association.Negative-Association-Permutation-Distributions*

**begin**

**unbundle** *no vec-syntax*

**context**

**fixes**  $f :: \text{real}$  **and**  $n :: \text{nat}$

**assumes** *f-range*:  $\langle f \in \{1/2..<1\} \rangle$   $\langle n * f \in \mathbb{N} \rangle$  **and** *n-gt-0*:  $\langle n > 0 \rangle$

**begin**

**definition**  $\langle \text{initial-state} = \text{State } \{\} \ 1 \rangle$  — Setup initial state  $\chi = \emptyset$  and  $p = 1$ .

**fun** *subsample* **where** — Subsampling operation: Sample random  $nf$  subset.

$\langle \text{subsample } \chi = \text{pmf-of-set } \{S. S \subseteq \chi \wedge \text{card } S = n * f\} \rangle$

**fun** *step* **where** — Loop body.

$\langle \text{step } a \ (\text{State } \chi \ p) = \text{do } \{$   
 $\quad b \leftarrow \text{bernoulli-pmf } p;$   
 $\quad \text{let } \chi = (\text{if } b \text{ then } \chi \cup \{a\} \text{ else } \chi - \{a\});$

$\quad \text{if } \text{card } \chi = n \text{ then do } \{$   
 $\quad \quad \chi \leftarrow \text{subsample } \chi;$   
 $\quad \quad \text{return-pmf } (\text{State } \chi \ (p * f))$   
 $\quad \quad \} \text{ else do } \{$   
 $\quad \quad \text{return-pmf } (\text{State } \chi \ p)$   
 $\quad \quad \}$   
 $\quad \}$   
 $\rangle$

**fun** *run-steps* **where** — Iterate loop over stream  $xs$ .

$\langle \text{run-steps } xs = \text{foldM-pmf } \text{step } xs \ \text{initial-state} \rangle$

**fun** *estimate* **where**

$\langle \text{estimate } (\text{State } \chi \ p) = \text{card } \chi / p \rangle$

**fun** *run-algo* **where** — Run algorithm and estimate.

$\langle \text{run-algo } xs = \text{map-pmf } \text{estimate } (\text{run-steps } xs) \rangle$

**definition**  $\langle \text{subsample-size} = (\text{THE } x. \text{real } x = n * f) \rangle$

**declare** *subsample.simps* [*simp del*]

**lemma** *subsample-size-eq*:

$\langle \text{real } \text{subsample-size} = n * f \rangle$

$\langle \text{proof} \rangle$

**lemma** *subsample-size*:

$\langle \text{subsample-size} < n \rangle \langle 2 * \text{subsample-size} \geq n \rangle$

$\langle \text{proof} \rangle$

**lemma** *subsample-finite-nonempty*:

**assumes**  $\langle \text{card } U = n \rangle$

**shows**

$\langle \{T. T \subseteq U \wedge \text{card } T = \text{subsample-size}\} \neq \{\} \rangle$  (**is**  $\langle ?C \neq \{\} \rangle$ )

$\langle \text{finite } \{T. T \subseteq U \wedge \text{card } T = \text{subsample-size}\} \rangle$

$\langle \text{subsample } U = \text{pmf-of-set } \{T. T \subseteq U \wedge \text{card } T = \text{subsample-size}\} \rangle$

$\langle \text{finite } (\text{set-pmf } (\text{subsample } U)) \rangle$

$\langle \text{proof} \rangle$

**lemma** *int-prod-subsample-eq-prod-int*:

**fixes**  $g :: \langle \text{bool} \Rightarrow \text{real} \rangle$

**assumes**  $\langle \text{card } U = n \rangle \langle S \subseteq U \rangle \langle \text{range } g \subseteq \{0..\} \rangle$

**shows**  $\langle (\int \omega. (\prod s \in S. g(s \in \omega))) \ \partial \text{subsample } U \leq (\prod s \in S. (\int \omega. g \ \omega \ \partial \text{bernoulli-pmf } f)) \rangle$

(**is**  $\langle ?L \leq ?R \rangle$ )

$\langle \text{proof} \rangle$



**schematic-goal** *step-n-def*:  $\langle \text{step } x \ \sigma = ?x \rangle$   
 $\langle \text{proof} \rangle$

**interpretation** *abs: cvm-algo-abstract n f subsample*  
**rewrites**  $\langle \text{abs.run-steps} = \text{run-steps} \rangle$  **and**  $\langle \text{abs.estimate} = \text{estimate} \rangle$   
 $\langle \text{proof} \rangle$

**theorem** *unbiasedness*:  $\langle \text{measure-pmf.expectation} (\text{run-algo } xs) \ \text{id} = \text{card} (\text{set } xs) \rangle$   
 $\langle \text{proof} \rangle$

**theorem** *correctness*:  
**assumes**  $\langle \varepsilon \in \{0 < .. < 1 :: \text{real}\} \rangle$   $\langle \delta \in \{0 < .. < 1 :: \text{real}\} \rangle$   
**assumes**  $\langle \text{real } n \geq 12 / \varepsilon^2 * \ln (3 * \text{real} (\text{length } xs) / \delta) \rangle$   
**defines**  $\langle A \equiv \text{real} (\text{card} (\text{set } xs)) \rangle$   
**shows**  $\langle \mathcal{P}(R \text{ in run-algo } xs. |R - A| > \varepsilon * A) \leq \delta \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *space-usage*:  
 $\langle \text{AE } \sigma \text{ in run-steps } xs. \text{card} (\text{state-}\chi \ \sigma) < n \wedge \text{finite} (\text{state-}\chi \ \sigma) \rangle$   
 $\langle \text{proof} \rangle$

**end**

**end**

## References

- [1] S. Chakraborty, N. V. Vinodchandran, and K. S. Meel. Distinct elements in streams: An algorithm for the (text) book. In S. Chechik, G. Navarro, E. Rotenberg, and G. Herman, editors, *ESA*, volume 244 of *LIPICs*, pages 34:1–34:6. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2022.
- [2] D. P. Dubhashi, V. Priebe, and D. Ranjan. Negative dependence through the fkg inequality. *BRICS Report Series*, 3, 1996.
- [3] E. Karayel. Negatively associated random variables. *Archive of Formal Proofs*, January 2025. [https://isa-afp.org/entries/Negative\\_Association.html](https://isa-afp.org/entries/Negative_Association.html), Formal proof development.

## A Informal Proof

This section includes an informal version of the proof for the tail bounds and unbiasedness of the abstract algorithm (Algorithm 1) for interested readers.

This means we assume the  $\text{subsample}(\chi)$  operation fulfills Eq. 1 and always returns a subset of  $\chi$ .

**Notation:** For a finite set  $S$ , the probability space of uniformly sampling from the set is denoted by  $U(S)$ , i.e., for each  $s \in S$  we have  $\mathcal{P}_{U(S)}(s) = |S|^{-1}$ . We write  $\text{Ber}(p)$  for the Bernoulli probability space, over the set  $\{0, 1\}$ , i.e.,  $P_{\text{Ber}(p)}(\{1\}) = p$ .  $I(P)$  is the indicator function for a predicate  $P$ , i.e.,  $I(\text{true}) = 1$  and  $I(\text{false}) = 0$ .

Like in the formalization, we will denote the first five lines of the loop (3–7) as step 1, the last four lines (8–10) as step 2. For the distribution of the state of the algorithm after processing  $i$  elements of the sequence, we will write  $\Omega_i$ . The elements of the probability spaces are pairs composed of a set and the number of subsampling steps, representing  $\chi$  and  $p$  respectively.

For example:  $\Omega_0 = U(\{(\emptyset, 1)\})$  is the initial state,  $\Omega_1 = U(\{(\{a_1\}, 1)\})$ , etc., and  $\Omega_l$  denotes the final state. We introduce  $\chi$  and  $p$  as random variables defined over such probability spaces  $\Omega$ , in particular,  $\chi$  (resp.  $p$ ) is the projection to the first (resp. second) component.

The state of the algorithm after processing only step 1 of the  $i$ -th loop iteration is denoted by  $\Omega'_i$ . So the sequence of states is represented by the distributions  $\Omega_0, \Omega'_1, \Omega_1, \dots, \Omega'_l, \Omega_l$ .

### A.1 Loop Invariant

After these preliminaries, we can go to the main proof, whose core is a probabilistic loop invariant for Algorithm 1 that can be verified inductively.

**Lemma 1.** Let  $\varphi : (0, 1] \times \{0, 1\} \rightarrow \mathbb{R}_{\geq 0}$  be a function, fulfilling the following conditions:

1.  $(1 - \alpha)\varphi(x, 0) + \alpha\varphi(x, 1) \leq \varphi(x/\alpha, 1)$  for all  $0 < \alpha < 1$ ,  $0 < x \leq 1$ , and
2.  $\varphi(x, 0) \leq \varphi(y, 0)$  for all  $0 < x < y \leq 1$ .

Then for all  $k \in \{0, \dots, l\}$ ,  $S \subseteq \{a_1, \dots, a_k\}$ ,  $\Omega \in \{\Omega_k, \Omega'_k\}$ :

$$\mathbb{E}_{\Omega} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] \leq \varphi(1, 1)^{|S|}$$

*Proof.* We show the result using induction over  $k$ . Note that we show the statement for arbitrary  $S$ , i.e., the induction statements are:

$$P(k) \quad :\Leftrightarrow \quad \left( \forall S \subseteq \{a_1, \dots, a_k\}. \mathbb{E}_{\Omega_k} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] \leq \varphi(1, 1)^{|S|} \right)$$

$$Q(k) \quad :\Leftrightarrow \quad \left( \forall S \subseteq \{a_1, \dots, a_k\}. \mathbb{E}_{\Omega'_k} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] \leq \varphi(1, 1)^{|S|} \right)$$

and we will show  $P(0), Q(1), P(1), Q(2), P(2), \dots, Q(l), P(l)$  successively.

**Induction start  $P(0)$ :**

We have  $S \subseteq \emptyset$ , and hence

$$\mathbb{E}_{\Omega_0} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] = \mathbb{E}_{\Omega_0} [1] = 1 \leq \varphi(1, 1)^0.$$

**Induction step  $P(k) \rightarrow Q(k+1)$ :**

Let  $S \subseteq \{a_1, \dots, a_{k+1}\}$  and define  $S' := S - \{a_{k+1}\}$ . Note that  $\Omega'_{k+1}$  can be constructed from  $\Omega_k$  as a compound distribution, where  $a_{k+1}$  is included in the buffer, with the probability  $p$ , which is itself a random variable over the space  $\Omega_k$ .

In particular, for example:

$$\mathcal{P}_{\Omega'_{k+1}}(P(\chi, p)) = \int_{\Omega_k} \int_{\text{Ber}(p(\omega))} P(\chi(\omega) - \{a_{k+1}\} \cup f(\tau), p(\omega)) d\tau d\omega$$

for all predicates  $P$  where we define  $f(1) = \{a_{k+1}\}$  and  $f(0) = \emptyset$ .

We distinguish the two cases  $a_{k+1} \in S$  and  $a_{k+1} \notin S$ . If  $a_{k+1} \in S$ :

$$\begin{aligned} & \mathbb{E}_{\Omega'_{k+1}} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] \\ &= \int_{\Omega_k} \left( \prod_{s \in S'} \varphi(p, I(s \in \chi)) \right) \int_{\text{Ber}(p(\omega))} \varphi(p, \tau) d\tau d\omega \\ &= \int_{\Omega_k} \left( \prod_{s \in S'} \varphi(p, I(s \in \chi)) \right) ((1-p)\varphi(p, 0) + p\varphi(p, 1)) d\omega \\ &\stackrel{\text{Cond 1}}{\leq} \int_{\Omega_k} \left( \prod_{s \in S'} \varphi(p, I(s \in \chi)) \right) \varphi(1, 1) d\omega \\ &\stackrel{\text{IH}}{\leq} \varphi(1, 1)^{|S'|} \varphi(1, 1) = \varphi(1, 1)^{|S|} \end{aligned}$$

If  $a_{k+1} \notin S$  then  $S' = S$  and:

$$\mathbb{E}_{\Omega'_{k+1}} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] = \int_{\Omega_k} \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega \stackrel{\text{IH}}{\leq} \varphi(1, 1)^{|S'|} = \varphi(1, 1)^{|S|}$$

**Induction step  $Q(k+1) \rightarrow P(k+1)$ :**

Let  $S \subseteq \{a_1, \dots, a_{k+1}\}$ .

Let us again note that  $\Omega_{k+1}$  is a compound distribution over  $\Omega'_{k+1}$ . In general, for all predicates  $P$ :

$$\mathcal{P}_{\Omega_{k+1}}(P(\chi, p)) = \int_{\Omega'_{k+1}} I(|\chi(\omega)| < n) P(\chi(\omega), p(\omega)) + I(|\chi(\omega)| = n) \int_{\text{subsample}(\chi(\omega))} P(\tau, fp(\omega)) d\tau d\omega.$$

With this we can now verify the induction step:

$$\begin{aligned}
& \mathbb{E}_{\Omega_{k+1}} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] \\
&= \int_{\Omega'_{k+1}} I(|\chi| < n) \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega \\
&+ \int_{\Omega'_{k+1}} I(|\chi| = n) \prod_{s \in S \setminus \chi(\omega)} \varphi(pf, 0) \int_{\text{subsample}(\chi)} \prod_{s \in S \cap \chi} \varphi(pf, I(s \in \tau)) d\tau d\omega \\
&\leq \int_{\Omega'_{k+1}} I(|\chi| < n) \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega && \text{Eq. 1} \\
&+ \int_{\Omega'_{k+1}} I(|\chi| = n) \prod_{s \in S \setminus \chi(\omega)} \varphi(pf, 0) \prod_{s \in S \cap \chi} \int_{\text{Ber}(f)} \varphi(pf, \tau) d\tau d\omega \\
&\leq \int_{\Omega'_{k+1}} I(|\chi| < n) \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega && \text{Cond 2} \\
&+ \int_{\Omega'_{k+1}} I(|\chi| = n) \prod_{s \in S \setminus \chi(\omega)} \varphi(p, 0) \prod_{s \in S \cap \chi} ((1-f)\varphi(pf, 0) + f\varphi(pf, 1)) d\omega \\
&\leq \int_{\Omega'_{k+1}} I(|\chi| < n) \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega && \text{Cond 1} \\
&+ \int_{\Omega'_{k+1}} I(|\chi| = n) \prod_{s \in S \setminus \chi(\omega)} \varphi(p, 0) \prod_{s \in S \cap \chi} \varphi(p, 1) d\omega \\
&= \int_{\Omega'_{k+1}} I(|\chi| < n) \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega \\
&+ \int_{\Omega'_{k+1}} I(|\chi| = n) \prod_{s \in S} \varphi(p, I(s \in \chi)) d\omega \\
&= \mathbb{E}_{\Omega'_{k+1}} \left[ \prod_{s \in S} \varphi(p, I(s \in \chi)) \right] \leq \varphi(1, 1)^{|S|} && \text{IH}
\end{aligned}$$

□

A corollary and more practical version of the previous lemma is:

**Lemma 2.** Let  $q \leq 1$  and  $h : [0, q^{-1}] \rightarrow \mathbb{R}_{\geq 0}$  be concave then for all  $k \in \{0, \dots, l\}$ ,  $S \subseteq \{a_1, \dots, a_k\}$ ,  $\Omega \in \{\Omega_k, \Omega'_k\}$ :

$$\mathbb{E}_{\Omega} \left[ \prod_{s \in S} I(p > q) h(p^{-1} I(s \in \chi)) \right] \leq h(1)^{|S|}$$

*Proof.* Follows from Lemma 1 for  $\varphi(p, \tau) := I(p > q) h(\tau p^{-1})$ . We just need to check Conditions 1 and 2. Indeed,

$$\begin{aligned}
(1 - \alpha)\varphi(x, 0) + \alpha\varphi(x, 1) &= (1 - \alpha)I(x > q)h(0) + \alpha I(x > q)h(x^{-1}) \\
&\leq I(x > q)h(\alpha x^{-1}) \leq I(x > q\alpha)h(\alpha x^{-1}) = \varphi(x/\alpha, 1)
\end{aligned}$$

for  $0 < \alpha < 1$  and  $0 < x \leq 1$ , where we used that  $x > q$  implies  $x > q\alpha$ ; and

$$\varphi(x, 0) = I(x > q)h(0) \leq I(y > q)h(0) = \varphi(y, 0)$$

for  $0 < x < y \leq 1$ , where we used that  $x > q$  implies  $y > q$ . □

It should be noted that this is a probabilistic recurrence relation, but the main innovation is that we establish a relation, with respect to general classes of functions of the state variables.

## A.2 Concentration

Let us now see how we can obtain concentration bounds using Lemma 2, i.e., that the result of the algorithm is concentrated around the cardinality of  $A = \{a_1, \dots, a_l\}$ . This will be done using the Cramér–Chernoff method for the probability that the estimate is above  $(1 + \varepsilon)|A|$  (resp. below  $(1 - \varepsilon)|A|$ ) assuming  $p$  is not too small and a tail estimate for  $p$  being too small.

It should be noted that concentration is trivial, if  $|A| < n$ , i.e., if we never need to do sub-sampling, so we assume  $|A| \geq n$ .

Define  $q := n/(4|A|)$  and notice that  $q \leq \frac{1}{4}$ .

Let us start with the upper tail bound:

**Lemma 3.** For any  $\Omega \in \{\Omega_0, \dots, \Omega_l\} \cup \{\Omega'_1, \dots, \Omega'_l\}$  and  $0 < \varepsilon \leq 1$ :

$$L := \mathcal{P}_\Omega (p^{-1}|\chi| \geq (1 + \varepsilon)|A| \wedge p \geq q) \leq \exp\left(-\frac{n}{12}\varepsilon^2\right)$$

*Proof.* By assumption there exists a  $k$  such that  $\Omega \in \{\Omega_k, \Omega'_k\}$ . Let  $A' = A \cap \{a_1, \dots, a_k\}$ . Moreover, we define:

$$\begin{aligned} t &:= q \ln(1 + \varepsilon) \\ h(x) &:= 1 + qx(e^{t/q} - 1) \end{aligned}$$

To get a tail estimate, we use the Cramér–Chernoff method:

$$\begin{aligned} L &\stackrel{t>0}{\leq} \mathcal{P}_\Omega (\exp(tp^{-1}|\chi|) \geq \exp(t(1 + \varepsilon)|A|) \wedge p \geq q) \\ &\leq \mathcal{P}_\Omega (I(p \geq q) \exp(tp^{-1}|\chi|) \geq \exp(t(1 + \varepsilon)|A|)) \\ &\stackrel{\text{Markov}}{\leq} \exp(-t(1 + \varepsilon)|A|) \mathbb{E}_\Omega [I(p \geq q) \exp(tp^{-1}|\chi|)] \\ &\leq \exp(-t(1 + \varepsilon)|A|) \mathbb{E}_\Omega \left[ \prod_{s \in A'} I(p \geq q) \exp(tp^{-1}I(s \in \chi)) \right] \\ &\leq \exp(-t(1 + \varepsilon)|A|) \mathbb{E}_\Omega \left[ \prod_{s \in A'} I(p \geq q) h(p^{-1}I(s \in \chi)) \right] \\ &\stackrel{\text{Lé. 2}}{\leq} \exp(-t(1 + \varepsilon)|A|) h(1)^{|A'|} \\ &\stackrel{h(1) \geq 1}{\leq} (\exp(\ln(h(1)) - t(1 + \varepsilon)))^{|A|} \end{aligned}$$

So we just need to show that (using  $|A| = \frac{n}{4q}$ ):

$$\ln(h(1)) - t(1 + \varepsilon) \leq \frac{-q\varepsilon^2}{3}$$

The latter can be established by analyzing the function

$$f(\varepsilon) := -\ln(1 + q\varepsilon) + q \ln(1 + \varepsilon)(1 + \varepsilon) - \frac{q\varepsilon^2}{3} = -\ln(h(1)) + t(1 + \varepsilon) - \frac{q\varepsilon^2}{3}.$$

For which it is easy to check  $f(0) = 0$  and the derivative with respect to  $\varepsilon$  is non-negative in the range  $0 \leq q \leq 1/4$  and  $0 < \varepsilon \leq 1$ , i.e.,  $f(\varepsilon) \geq 0$ .  $\square$

Using the previous result we can also estimate bounds for  $p$  becoming too small:

**Lemma 4.**

$$\mathcal{P}_{\Omega_l}(p < q) \leq l \exp\left(-\frac{n}{12}\right)$$

*Proof.* We will use a similar strategy as in the  $\text{Bad}_2$  bound from the original CVM paper [1]. Let  $j$  be maximal, s.t.,  $q \leq f^j$ . Hence  $f^{j+1} < q$  and:

$$f^j \leq 2ff^j < 2q = \frac{n}{2|A|}. \quad (2)$$

First, we bound the probability of jumping from  $p = f^j$  to  $p = f^{j+1}$  at a specific point in the algorithm, e.g., while processing  $k$  stream elements. It can only happen if  $|\chi| = n$ ,  $p = f^j$  in  $\Omega'_k$ . Then

$$\begin{aligned} \mathcal{P}_{\Omega'_k}(|\chi| \geq n \wedge p = f^j) &\leq \mathcal{P}(p^{-1}|\chi| \geq f^{-j}n \wedge p \geq q) \\ &\stackrel{\text{Eq. 2}}{\leq} \mathcal{P}(p^{-1}|\chi| \geq 2|A| \wedge p \geq q) \\ &\stackrel{\text{Le. 3}}{\leq} \exp(-n/12) \end{aligned}$$

The probability that this happens ever in the entire process is then at most  $l$  times the above which proves the lemma.  $\square$

**Lemma 5.** Let  $0 < \varepsilon < 1$  then:

$$L := \mathcal{P}_{\Omega_l}(p^{-1}|\chi| \leq (1 - \varepsilon)|A| \wedge p \geq q) \leq \exp\left(-\frac{n}{8}\varepsilon^2\right)$$

*Proof.* Let us define

$$\begin{aligned} t &:= q \ln(1 - \varepsilon) < 0 \\ h(x) &:= 1 + qx(e^{t/q} - 1) \end{aligned}$$

Note that  $h(x) \geq 0$  for  $0 \leq x \leq q^{-1}$  (can be checked by verifying it is true for  $h(0)$  and  $h(q^{-1})$  and the fact that the function is affine.)

With these definitions we again follow the Cramér–Chernoff method:

$$\begin{aligned} L &\stackrel{t < 0}{=} \mathcal{P}_{\Omega_l}(\exp(tp^{-1}|\chi|) \geq \exp(t(1 - \varepsilon)|A|) \wedge p \geq q) \\ &\leq \mathcal{P}_{\Omega_l}(I(p \geq q) \exp(tp^{-1}|\chi|) \geq \exp(t(1 - \varepsilon)|A|) \wedge p > q) \\ &\stackrel{\text{Markov}}{\leq} \exp(-t(1 - \varepsilon)|A|) \mathbb{E}_{\Omega} [I(p \geq q) \exp(tp^{-1}|\chi|)] \\ &= \exp(-t(1 - \varepsilon)|A|) \mathbb{E}_{\Omega} \left[ \prod_{s \in A} I(p \geq q) \exp(tp^{-1}I(s \in \chi)) \right] \\ &\leq \exp(-t(1 - \varepsilon)|A|) \mathbb{E}_{\Omega} \left[ \prod_{s \in A} I(p \geq q) h(p^{-1}I(s \in \chi)) \right] \\ &\stackrel{\text{Le. 2}}{\leq} \exp(-t(1 - \varepsilon)|A|) (h(1))^{|A|} \\ &= \exp(\ln(h(1)) - t(1 - \varepsilon))^{|A|} \end{aligned}$$

Substituting  $t$  and  $h$  and using  $|A| = \frac{n}{4q}$ , we can see that the lemma is true if

$$f(\varepsilon) := q \ln(1 - \varepsilon)(1 - \varepsilon) - \ln(1 - q\varepsilon) - \frac{q}{2}\varepsilon^2 = t(1 - \varepsilon) - \ln(h(1)) - \frac{q}{2}\varepsilon^2$$

is non-negative for  $0 \leq q \leq \frac{1}{4}$  and  $0 < \varepsilon < 1$ . This can be verified by checking that  $f(0) = 0$  and that the derivative with respect to  $\varepsilon$  is non-negative.  $\square$

We can now establish the concentration result:

**Theorem 1.** *Let  $0 < \varepsilon < 1$  and  $0 < \delta < 1$  and  $n \geq \frac{12}{\varepsilon^2} \ln\left(\frac{3l}{\delta}\right)$  then:*

$$L = \mathcal{P}_{\Omega_l}(|p^{-1}|\chi| - |A| \geq \varepsilon|A|) \leq \delta$$

*Proof.* Note that the theorem is trivial if  $|A| < n$ . If not:

$$\begin{aligned} L &\leq \mathcal{P}_{\Omega_l}(|p^{-1}|\chi| \leq (1 - \varepsilon)|A| \wedge p \geq q) \\ &\quad + \mathcal{P}_{\Omega_l}(|p^{-1}|\chi| \geq (1 + \varepsilon)|A| \wedge p \geq q) + \mathcal{P}_{\Omega_l}(p < q) \\ &\stackrel{\text{L.e. 3-5}}{\leq} \exp\left(-\frac{n}{8}\varepsilon^2\right) + \exp\left(-\frac{n}{12}\varepsilon^2\right) + l \exp\left(-\frac{n}{12}\right) \\ &\leq \frac{\delta}{3} + \frac{\delta}{3} + \frac{\delta}{3} \end{aligned}$$

$\square$

### A.3 Unbiasedness

Let  $M$  be large enough such that  $p^{-1} \leq M$  a.s. (e.g., we can choose  $M = f^{-l}$ ). Then we can derive from Lemma 2 using  $h(x) = x$  and  $h(x) = M + 1 - x$  that for all  $s \in A$ :

$$\begin{aligned} \mathbb{E}_{\Omega_l}[p^{-1}I(s \in \chi)] &= \mathbb{E}_{\Omega_l}[I(p \geq M^{-1})p^{-1}I(s \in \chi)] \leq 1 \\ \mathbb{E}_{\Omega_l}[M + 1 - p^{-1}I(s \in \chi)] &= \mathbb{E}_{\Omega_l}[I(p \geq M^{-1})(M + 1 - p^{-1}I(s \in \chi))] \leq M \end{aligned}$$

which implies  $\mathbb{E}_{\Omega_l}[p^{-1}I(s \in \chi)] = 1$ . By linearity of expectation we conclude

$$\mathbb{E}_{\Omega_l}[p^{-1}|\chi|] = \sum_{s \in A} \mathbb{E}_{\Omega_l}[p^{-1}I(s \in \chi)] = |A|.$$