

# Formalization of Randomized Approximation Algorithms for Frequency Moments

Emin Karayel

March 11, 2024

## Abstract

In 1999 Alon et. al. introduced the still active research topic of approximating the frequency moments of a data stream using randomized algorithms with minimal space usage. This includes the problem of estimating the cardinality of the stream elements—the zeroth frequency moment. But, also higher-order frequency moments that provide information about the skew of the data stream. (The  $k$ -th frequency moment of a data stream is the sum of the  $k$ -th powers of the occurrence counts of each element in the stream.) This entry formalizes three randomized algorithms for the approximation of  $F_0$ ,  $F_2$  and  $F_k$  for  $k \geq 3$  based on [1, 2] and verifies their expected accuracy, success probability and space usage.

## Contents

<b>1</b>	<b>Preliminary Results</b>	<b>2</b>
<b>2</b>	<b>Frequency Moments</b>	<b>5</b>
<b>3</b>	<b>Ranks, <math>k</math> smallest element and elements</b>	<b>6</b>
<b>4</b>	<b>Landau Symbols</b>	<b>8</b>
<b>5</b>	<b>Probability Spaces</b>	<b>10</b>
<b>6</b>	<b>Indexed Products of Probability Mass Functions</b>	<b>12</b>
<b>7</b>	<b>Frequency Moment 0</b>	<b>12</b>
<b>8</b>	<b>Frequency Moment 2</b>	<b>16</b>
<b>9</b>	<b>Frequency Moment <math>k</math></b>	<b>21</b>

<b>A Informal proof of correctness for the <math>F_0</math> algorithm</b>	<b>26</b>
A.1 Case $F_0 \geq t$ . . . . .	27
A.2 Case $F_0 < t$ . . . . .	29

## 1 Preliminary Results

**theory** *Frequency-Moments-Preliminary-Results*

**imports**

*HOL.Transcendental*

*HOL-Computational-Algebra.Primes*

*HOL-Library.Extended-Real*

*HOL-Library.Multiset*

*HOL-Library.Sublist*

*Prefix-Free-Code-Combinators.Prefix-Free-Code-Combinators*

*Bertrands-Postulate.Bertrand*

*Expander-Graphs.Expander-Graphs-Multiset-Extras*

**begin**

This section contains various preliminary results.

**lemma** *card-ordered-pairs*:

**fixes**  $M :: ('a :: \text{linorder}) \text{ set}$

**assumes** *finite M*

**shows**  $2 * \text{card } \{(x,y) \in M \times M. x < y\} = \text{card } M * (\text{card } M - 1)$

*<proof>*

**lemma** *ereal-mono*:  $x \leq y \implies \text{ereal } x \leq \text{ereal } y$

*<proof>*

**lemma** *log-mono*:  $a > 1 \implies x \leq y \implies 0 < x \implies \log a \ x \leq \log a \ y$

*<proof>*

**lemma** *abs-ge-iff*:  $((x::\text{real}) \leq \text{abs } y) = (x \leq y \vee x \leq -y)$

*<proof>*

**lemma** *count-list-gr-1*:

$(x \in \text{set } xs) = (\text{count-list } xs \ x \geq 1)$

*<proof>*

**lemma** *count-list-append*:  $\text{count-list } (xs@ys) \ v = \text{count-list } xs \ v + \text{count-list } ys \ v$

*<proof>*

**lemma** *count-list-lt-suffix*:

**assumes** *suffix a b*

**assumes**  $x \in \{b \ ! \ i \mid i. i < \text{length } b - \text{length } a\}$

**shows**  $\text{count-list } a \ x < \text{count-list } b \ x$

*<proof>*

**lemma** *suffix-drop-drop*:

**assumes**  $x \geq y$   
**shows**  $\text{suffix } (\text{drop } x \ a) \ (\text{drop } y \ a)$   
 $\langle \text{proof} \rangle$

**lemma** *count-list-card*:  $\text{count-list } xs \ x = \text{card } \{k. k < \text{length } xs \wedge xs ! k = x\}$   
 $\langle \text{proof} \rangle$

**lemma** *card-gr-1-iff*:  
**assumes**  $\text{finite } S \ x \in S \ y \in S \ x \neq y$   
**shows**  $\text{card } S > 1$   
 $\langle \text{proof} \rangle$

**lemma** *count-list-ge-2-iff*:  
**assumes**  $y < z$   
**assumes**  $z < \text{length } xs$   
**assumes**  $xs ! y = xs ! z$   
**shows**  $\text{count-list } xs \ (xs ! y) > 1$   
 $\langle \text{proof} \rangle$

Results about multisets and sorting

**lemmas** *disj-induct-mset* = *disj-induct-mset*

**lemma** *prod-mset-conv*:  
**fixes**  $f :: 'a \Rightarrow 'b::\{\text{comm-monoid-mult}\}$   
**shows**  $\text{prod-mset } (\text{image-mset } f \ A) = \text{prod } (\lambda x. f \ x \ (\text{count } A \ x)) \ (\text{set-mset } A)$   
 $\langle \text{proof} \rangle$

There is a version *sum-list-map-eq-sum-count* but it doesn't work if the function maps into the reals.

**lemma** *sum-list-eval*:  
**fixes**  $f :: 'a \Rightarrow 'b::\{\text{ring}, \text{semiring-1}\}$   
**shows**  $\text{sum-list } (\text{map } f \ xs) = (\sum x \in \text{set } xs. \text{of-nat } (\text{count-list } xs \ x) * f \ x)$   
 $\langle \text{proof} \rangle$

**lemma** *prod-list-eval*:  
**fixes**  $f :: 'a \Rightarrow 'b::\{\text{ring}, \text{semiring-1}, \text{comm-monoid-mult}\}$   
**shows**  $\text{prod-list } (\text{map } f \ xs) = (\prod x \in \text{set } xs. (f \ x) \ (\text{count-list } xs \ x))$   
 $\langle \text{proof} \rangle$

**lemma** *sorted-sorted-list-of-multiset*:  $\text{sorted } (\text{sorted-list-of-multiset } M)$   
 $\langle \text{proof} \rangle$

**lemma** *count-mset*:  $\text{count } (\text{mset } xs) \ a = \text{count-list } xs \ a$   
 $\langle \text{proof} \rangle$

**lemma** *swap-filter-image*:  $\text{filter-mset } g \ (\text{image-mset } f \ A) = \text{image-mset } f \ (\text{filter-mset } (g \circ f) \ A)$   
 $\langle \text{proof} \rangle$

**lemma** *list-eq-iff*:

**assumes** *mset xs = mset ys*

**assumes** *sorted xs*

**assumes** *sorted ys*

**shows** *xs = ys*

*<proof>*

**lemma** *sorted-list-of-multiset-image-commute*:

**assumes** *mono f*

**shows** *sorted-list-of-multiset (image-mset f M) = map f (sorted-list-of-multiset M)*

*<proof>*

Results about rounding and floating point numbers

**lemma** *round-down-ge*:

$x \leq \text{round-down } \text{prec } x + 2^{\text{powr } (-\text{prec})}$

*<proof>*

**lemma** *truncate-down-ge*:

$x \leq \text{truncate-down } \text{prec } x + \text{abs } x * 2^{\text{powr } (-\text{prec})}$

*<proof>*

**lemma** *truncate-down-pos*:

**assumes**  $x \geq 0$

**shows**  $x * (1 - 2^{\text{powr } (-\text{prec})}) \leq \text{truncate-down } \text{prec } x$

*<proof>*

**lemma** *truncate-down-eq*:

**assumes** *truncate-down r x = truncate-down r y*

**shows**  $\text{abs } (x - y) \leq \max (\text{abs } x) (\text{abs } y) * 2^{\text{powr } (-\text{real } r)}$

*<proof>*

**definition** *rat-of-float* :: *float*  $\Rightarrow$  *rat* **where**

*rat-of-float f = of-int (mantissa f) \**

*(if exponent f  $\geq 0$  then  $2^{\text{nat } (\text{exponent } f)}$  else  $1 / 2^{\text{nat } (-\text{exponent } f)}$ ))*

**lemma** *real-of-rat-of-float*: *real-of-rat (rat-of-float x) = real-of-float x*

*<proof>*

**lemma** *log-est*:  $\log 2 (\text{real } n + 1) \leq n$

*<proof>*

**lemma** *truncate-mantissa-bound*:

$\text{abs } (\lfloor x * 2^{\text{powr } (\text{real } r - \text{real-of-int } \lfloor \log 2 |x| \rfloor)} \rfloor) \leq 2^{\wedge(r+1)}$  (**is** ?lhs  $\leq$  -)

*<proof>*

**lemma** *truncate-float-bit-count*:

$\text{bit-count } (F_e (\text{float-of } (\text{truncate-down } r x))) \leq 10 + 4 * \text{real } r + 2 * \log 2 (2 +$

```

|log 2 |x||)
(is ?lhs ≤ ?rhs)
⟨proof⟩

```

```

definition prime-above :: nat ⇒ nat
  where prime-above n = (SOME x. x ∈ {n..(2*n+2)} ∧ prime x)

```

The term *prime-above*  $n$  returns a prime between  $n$  and  $2 * n + 2$ . Because of Bertrand's postulate there always is such a value. In a refinement of the algorithms, it may make sense to replace this with an algorithm, that finds such a prime exactly or approximately.

The definition is intentionally inexact, to allow refinement with various algorithms, without modifying the high-level mathematical correctness proof.

```

lemma ex-subset:
  assumes ∃ x ∈ A. P x
  assumes A ⊆ B
  shows ∃ x ∈ B. P x
  ⟨proof⟩

```

```

lemma
  shows prime-above-prime: prime (prime-above n)
  and prime-above-range: prime-above n ∈ {n..(2*n+2)}
  ⟨proof⟩

```

```

lemma prime-above-min: prime-above n ≥ 2
  ⟨proof⟩

```

```

lemma prime-above-lower-bound: prime-above n ≥ n
  ⟨proof⟩

```

```

lemma prime-above-upper-bound: prime-above n ≤ 2*n+2
  ⟨proof⟩

```

```

end

```

## 2 Frequency Moments

```

theory Frequency-Moments
  imports
    Frequency-Moments-Preliminary-Results
    Universal-Hash-Families.Universal-Hash-Families-More-Finite-Fields
    Interpolation-Polynomials-HOL-Algebra.Interpolation-Polynomial-Cardinalities
  begin

```

This section contains a definition of the frequency moments of a stream and a few general results about frequency moments..

```

definition F where

```

$$F\ k\ xs = (\sum\ x \in\ set\ xs.\ (rat-of-nat\ (count-list\ xs\ x)\ k))$$

**lemma** *F-ge-0*:  $F\ k\ as \geq 0$   
 $\langle proof \rangle$

**lemma** *F-gr-0*:  
**assumes**  $as \neq []$   
**shows**  $F\ k\ as > 0$   
 $\langle proof \rangle$

**definition**  $P_e :: nat \Rightarrow nat \Rightarrow nat\ list \Rightarrow bool\ list\ option$  **where**  
 $P_e\ p\ n\ f = (if\ p > 1 \wedge f \in bounded-degree-polynomials\ (mod-ring\ p)\ n\ then$   
 $([0..<n] \rightarrow_e Nb_e\ p)\ (\lambda i \in \{..<n\}.\ ring.coeff\ (mod-ring\ p)\ f\ i)\ else\ None)$

**lemma** *poly-encoding*:  
 $is-encoding\ (P_e\ p\ n)$   
 $\langle proof \rangle$

**lemma** *bounded-degree-polynomial-bit-count*:  
**assumes**  $p > 1$   
**assumes**  $x \in bounded-degree-polynomials\ (mod-ring\ p)\ n$   
**shows**  $bit-count\ (P_e\ p\ n\ x) \leq ereal\ (real\ n * (\log\ 2\ p + 1))$   
 $\langle proof \rangle$

**end**

### 3 Ranks, $k$ smallest element and elements

**theory** *K-Smallest*  
**imports**  
*Frequency-Moments-Preliminary-Results*  
*Interpolation-Polynomials-HOL-Algebra.Interpolation-Polynomial-Cardinalities*  
**begin**

This section contains definitions and results for the selection of the  $k$  smallest elements, the  $k$ -th smallest element, rank of an element in an ordered set.

**definition** *rank-of* ::  $'a :: linorder \Rightarrow 'a\ set \Rightarrow nat$  **where**  $rank-of\ x\ S = card\ \{y \in S.\ y < x\}$

The function *rank-of* returns the rank of an element within a set.

**lemma** *rank-mono*:  
**assumes**  $finite\ S$   
**shows**  $x \leq y \implies rank-of\ x\ S \leq rank-of\ y\ S$   
 $\langle proof \rangle$

**lemma** *rank-mono-2*:  
**assumes**  $finite\ S$   
**shows**  $S' \subseteq S \implies rank-of\ x\ S' \leq rank-of\ x\ S$

*<proof>*

**lemma** *rank-mono-commute*:

**assumes** *finite S*

**assumes**  $S \subseteq T$

**assumes** *strict-mono-on T f*

**assumes**  $x \in T$

**shows**  $\text{rank-of } x \ S = \text{rank-of } (f \ x) \ (f \ ' \ S)$

*<proof>*

**definition** *least* **where**  $\text{least } k \ S = \{y \in S. \text{rank-of } y \ S < k\}$

The function *K-Smallest.least* returns the k smallest elements of a finite set.

**lemma** *rank-strict-mono*:

**assumes** *finite S*

**shows** *strict-mono-on S*  $(\lambda x. \text{rank-of } x \ S)$

*<proof>*

**lemma** *rank-of-image*:

**assumes** *finite S*

**shows**  $(\lambda x. \text{rank-of } x \ S) \ ' \ S = \{0..<\text{card } S\}$

*<proof>*

**lemma** *card-least*:

**assumes** *finite S*

**shows**  $\text{card } (\text{least } k \ S) = \min k \ (\text{card } S)$

*<proof>*

**lemma** *least-subset*:  $\text{least } k \ S \subseteq S$

*<proof>*

**lemma** *least-mono-commute*:

**assumes** *finite S*

**assumes** *strict-mono-on S f*

**shows**  $f \ ' \ \text{least } k \ S = \text{least } k \ (f \ ' \ S)$

*<proof>*

**lemma** *least-eq-iff*:

**assumes** *finite B*

**assumes**  $A \subseteq B$

**assumes**  $\bigwedge x. x \in B \implies \text{rank-of } x \ B < k \implies x \in A$

**shows**  $\text{least } k \ A = \text{least } k \ B$

*<proof>*

**lemma** *least-insert*:

**assumes** *finite S*

**shows**  $\text{least } k \ (\text{insert } x \ (\text{least } k \ S)) = \text{least } k \ (\text{insert } x \ S)$  (**is** *?lhs = ?rhs*)

*<proof>*

**definition** *count-le* **where** *count-le*  $x\ M = \text{size } \{\#y \in \# M. y \leq x\# \}$   
**definition** *count-less* **where** *count-less*  $x\ M = \text{size } \{\#y \in \# M. y < x\# \}$

**definition** *nth-mset*  $:: \text{nat} \Rightarrow ('a :: \text{linorder}) \text{multiset} \Rightarrow 'a$  **where**  
*nth-mset*  $k\ M = \text{sorted-list-of-multiset } M\ !\ k$

**lemma** *nth-mset-bound-left*:  
**assumes**  $k < \text{size } M$   
**assumes** *count-less*  $x\ M \leq k$   
**shows**  $x \leq \text{nth-mset } k\ M$   
 $\langle \text{proof} \rangle$

**lemma** *nth-mset-bound-left-excl*:  
**assumes**  $k < \text{size } M$   
**assumes** *count-le*  $x\ M \leq k$   
**shows**  $x < \text{nth-mset } k\ M$   
 $\langle \text{proof} \rangle$

**lemma** *nth-mset-bound-right*:  
**assumes**  $k < \text{size } M$   
**assumes** *count-le*  $x\ M > k$   
**shows**  $\text{nth-mset } k\ M \leq x$   
 $\langle \text{proof} \rangle$

**lemma** *nth-mset-commute-mono*:  
**assumes** *mono*  $f$   
**assumes**  $k < \text{size } M$   
**shows**  $f\ (\text{nth-mset } k\ M) = \text{nth-mset } k\ (\text{image-mset } f\ M)$   
 $\langle \text{proof} \rangle$

**lemma** *nth-mset-max*:  
**assumes**  $\text{size } A > k$   
**assumes**  $\bigwedge x. x \leq \text{nth-mset } k\ A \implies \text{count } A\ x \leq 1$   
**shows**  $\text{nth-mset } k\ A = \text{Max } (\text{least } (k+1)\ (\text{set-mset } A))$  **and**  $\text{card } (\text{least } (k+1)\ (\text{set-mset } A)) = k+1$   
 $\langle \text{proof} \rangle$

**end**

## 4 Landau Symbols

**theory** *Landau-Ext*  
**imports**  
*HOL-Library.Landau-Symbols*  
*HOL.Topological-Spaces*  
**begin**

This section contains results about Landau Symbols in addition to "HOL-



Library.Landau".

**lemma** *landau-sum*:

**assumes** *eventually* ( $\lambda x. g1\ x \geq (0::real)$ ) *F*  
**assumes** *eventually* ( $\lambda x. g2\ x \geq 0$ ) *F*  
**assumes**  $f1 \in O[F](g1)$   
**assumes**  $f2 \in O[F](g2)$   
**shows** ( $\lambda x. f1\ x + f2\ x$ )  $\in O[F](\lambda x. g1\ x + g2\ x)$   
*<proof>*

**lemma** *landau-sum-1*:

**assumes** *eventually* ( $\lambda x. g1\ x \geq (0::real)$ ) *F*  
**assumes** *eventually* ( $\lambda x. g2\ x \geq 0$ ) *F*  
**assumes**  $f \in O[F](g1)$   
**shows**  $f \in O[F](\lambda x. g1\ x + g2\ x)$   
*<proof>*

**lemma** *landau-sum-2*:

**assumes** *eventually* ( $\lambda x. g1\ x \geq (0::real)$ ) *F*  
**assumes** *eventually* ( $\lambda x. g2\ x \geq 0$ ) *F*  
**assumes**  $f \in O[F](g2)$   
**shows**  $f \in O[F](\lambda x. g1\ x + g2\ x)$   
*<proof>*

**lemma** *landau-ln-3*:

**assumes** *eventually* ( $\lambda x. (1::real) \leq f\ x$ ) *F*  
**assumes**  $f \in O[F](g)$   
**shows** ( $\lambda x. \ln\ (f\ x)$ )  $\in O[F](g)$   
*<proof>*

**lemma** *landau-ln-2*:

**assumes**  $a > (1::real)$   
**assumes** *eventually* ( $\lambda x. 1 \leq f\ x$ ) *F*  
**assumes** *eventually* ( $\lambda x. a \leq g\ x$ ) *F*  
**assumes**  $f \in O[F](g)$   
**shows** ( $\lambda x. \ln\ (f\ x)$ )  $\in O[F](\lambda x. \ln\ (g\ x))$   
*<proof>*

**lemma** *landau-real-nat*:

**fixes**  $f :: 'a \Rightarrow int$   
**assumes** ( $\lambda x. of\_int\ (f\ x)$ )  $\in O[F](g)$   
**shows** ( $\lambda x. real\ (nat\ (f\ x))$ )  $\in O[F](g)$   
*<proof>*

**lemma** *landau-ceil*:

**assumes** ( $\lambda x. 1$ )  $\in O[F](g)$   
**assumes**  $f \in O[F](g)$   
**shows** ( $\lambda x. real\_of\_int\ \lceil f\ x \rceil$ )  $\in O[F](g)$   
*<proof>*

```

lemma landau-rat-ceil:
  assumes  $(\lambda -. 1) \in O[F^\uparrow](g)$ 
  assumes  $(\lambda x. \text{real-of-rat } (f\ x)) \in O[F^\uparrow](g)$ 
  shows  $(\lambda x. \text{real-of-int } \lceil f\ x \rceil) \in O[F^\uparrow](g)$ 
   $\langle \text{proof} \rangle$ 

lemma landau-nat-ceil:
  assumes  $(\lambda -. 1) \in O[F^\uparrow](g)$ 
  assumes  $f \in O[F^\uparrow](g)$ 
  shows  $(\lambda x. \text{real } (\text{nat } \lceil f\ x \rceil)) \in O[F^\uparrow](g)$ 
   $\langle \text{proof} \rangle$ 

lemma eventually-prod1':
  assumes  $B \neq \text{bot}$ 
  assumes  $(\forall_F x \text{ in } A. P\ x)$ 
  shows  $(\forall_F x \text{ in } A \times_F B. P\ (\text{fst } x))$ 
   $\langle \text{proof} \rangle$ 

lemma eventually-prod2':
  assumes  $A \neq \text{bot}$ 
  assumes  $(\forall_F x \text{ in } B. P\ x)$ 
  shows  $(\forall_F x \text{ in } A \times_F B. P\ (\text{snd } x))$ 
   $\langle \text{proof} \rangle$ 

lemma sequentially-inf:  $\forall_F x \text{ in sequentially. } n \leq \text{real } x$ 
   $\langle \text{proof} \rangle$ 

instantiation rat :: linorder-topology
begin

definition open-rat :: rat set  $\Rightarrow$  bool
  where open-rat = generate-topology (range  $(\lambda a. \{.. < a\}) \cup \text{range } (\lambda a. \{a <..\})$ )

instance
   $\langle \text{proof} \rangle$ 
end

lemma inv-at-right-0-inf:
   $\forall_F x \text{ in at-right } 0. c \leq 1 \text{ / real-of-rat } x$ 
   $\langle \text{proof} \rangle$ 

end

```

## 5 Probability Spaces

Some additional results about probability spaces in addition to "HOL-Probability".

```

theory Probability-Ext
  imports

```

*HOL-Probability.Stream-Space*  
*Concentration-Inequalities.Bienaymes-Identity*  
*Universal-Hash-Families.Carter-Wegman-Hash-Family*  
*Frequency-Moments-Preliminary-Results*

**begin**

The following aliases are here to prevent possible merge-conflicts. The lemmas have been moved to *Concentration-Inequalities.Bienaymes-Identity* and/or *Concentration-Inequalities.Concentration-Inequalities-Preliminary*.

**lemmas** *make-ext = forall-Pi-to-PiE*

**lemmas** *PiE-reindex = PiE-reindex*

**context** *prob-space*

**begin**

**lemmas** *indep-sets-reindex = indep-sets-reindex*

**lemmas** *indep-vars-cong-AE = indep-vars-cong-AE*

**lemmas** *indep-vars-reindex = indep-vars-reindex*

**lemmas** *variance-divide = variance-divide*

**lemmas** *covariance-def = covariance-def*

**lemmas** *real-prod-integrable = cauchy-schwartz(1)*

**lemmas** *covariance-eq = covariance-eq*

**lemmas** *covar-integrable = covar-integrable*

**lemmas** *sum-square-int = sum-square-int*

**lemmas** *var-sum-1 = bienaymes-identity*

**lemmas** *covar-self-eq = covar-self-eq*

**lemmas** *covar-indep-eq-zero = covar-indep-eq-zero*

**lemmas** *var-sum-2 = bienaymes-identity-2*

**lemmas** *var-sum-pairwise-indep = bienaymes-identity-pairwise-indep*

**lemmas** *indep-var-from-indep-vars = indep-var-from-indep-vars*

**lemmas** *var-sum-pairwise-indep-2 = bienaymes-identity-pairwise-indep-2*

**lemmas** *var-sum-all-indep = bienaymes-identity-full-indep*

**lemma** *pmf-mono:*

**assumes**  $M = \text{measure-pmf } p$

**assumes**  $\bigwedge x. x \in P \implies x \in \text{set-pmf } p \implies x \in Q$

**shows**  $\text{prob } P \leq \text{prob } Q$

*<proof>*

**lemma** *pmf-add:*

**assumes**  $M = \text{measure-pmf } p$

**assumes**  $\bigwedge x. x \in P \implies x \in \text{set-pmf } p \implies x \in Q \vee x \in R$

**shows**  $\text{prob } P \leq \text{prob } Q + \text{prob } R$

*<proof>*

**lemma** *pmf-add-2:*

**assumes**  $M = \text{measure-pmf } p$

**assumes**  $\text{prob } \{\omega. P \ \omega\} \leq r1$

**assumes**  $\text{prob } \{\omega. Q \ \omega\} \leq r2$

**shows**  $\text{prob } \{\omega. P \ \omega \vee Q \ \omega\} \leq r1 + r2$  (**is**  $?lhs \leq ?rhs$ )  
 $\langle \text{proof} \rangle$

**end**

**end**

## 6 Indexed Products of Probability Mass Functions

**theory** *Product-PMF-Ext*

**imports**

*Probability-Ext*

*Universal-Hash-Families.Universal-Hash-Families-More-Product-PMF*

**begin**

The following aliases are here to prevent possible merge-conflicts. The lemmas have been moved to *Universal-Hash-Families.Universal-Hash-Families-More-Product-PMF*.

**abbreviation** *prod-pmf* **where** *prod-pmf*  $\equiv$  *Universal-Hash-Families-More-Product-PMF.prod-pmf*

**abbreviation** *restrict-dfl* **where** *restrict-dfl*  $\equiv$  *Universal-Hash-Families-More-Product-PMF.restrict-dfl*

**lemmas** *pmf-prod-pmf* = *pmf-prod-pmf*

**lemmas** *PiE-default-undefined-eq* = *PiE-default-undefined-eq*

**lemmas** *set-prod-pmf* = *set-prod-pmf*

**lemmas** *prob-prod-pmf'* = *prob-prod-pmf'*

**lemmas** *prob-prod-pmf-slice* = *prob-prod-pmf-slice*

**lemmas** *pi-pmf-decompose* = *pi-pmf-decompose*

**lemmas** *restrict-dfl-iter* = *restrict-dfl-iter*

**lemmas** *indep-vars-restrict'* = *indep-vars-restrict'*

**lemmas** *indep-vars-restrict-intro'* = *indep-vars-restrict-intro'*

**lemmas** *integrable-Pi-pmf-slice* = *integrable-Pi-pmf-slice*

**lemmas** *expectation-Pi-pmf-slice* = *expectation-Pi-pmf-slice*

**lemmas** *expectation-prod-Pi-pmf* = *expectation-prod-Pi-pmf*

**lemmas** *variance-prod-pmf-slice* = *variance-prod-pmf-slice*

**lemmas** *Pi-pmf-bind-return* = *Pi-pmf-bind-return*

**end**

## 7 Frequency Moment 0

**theory** *Frequency-Moment-0*

**imports**

*Frequency-Moments-Preliminary-Results*

*Median-Method.Median*

*K-Smallest*

*Universal-Hash-Families.Carter-Wegman-Hash-Family*

*Frequency-Moments*

*Landau-Ext*

*Probability-Ext*

**begin**

This section contains a formalization of a new algorithm for the zero-th frequency moment inspired by ideas described in [2]. It is a KMV-type ( $k$ -minimum value) algorithm with a rounding method and matches the space complexity of the best algorithm described in [2].

In addition to the Isabelle proof here, there is also an informal hand-written proof in Appendix A.

**type-synonym**  $f0\text{-state} = \text{nat} \times \text{nat} \times \text{nat} \times \text{nat} \times (\text{nat} \Rightarrow \text{nat list}) \times (\text{nat} \Rightarrow \text{float set})$

**definition**  $\text{hash}$  **where**  $\text{hash } p = \text{ring.hash } (\text{mod-ring } p)$

**fun**  $f0\text{-init} :: \text{rat} \Rightarrow \text{rat} \Rightarrow \text{nat} \Rightarrow f0\text{-state pmf}$  **where**  
 $f0\text{-init } \delta \ \varepsilon \ n =$   
 do {  
   let  $s = \text{nat } \lceil -18 * \ln (\text{real-of-rat } \varepsilon) \rceil$ ;  
   let  $t = \text{nat } \lceil 80 / (\text{real-of-rat } \delta)^2 \rceil$ ;  
   let  $p = \text{prime-above } (\text{max } n \ 19)$ ;  
   let  $r = \text{nat } (4 * \lceil \log 2 (1 / \text{real-of-rat } \delta) \rceil + 23)$ ;  
    $h \leftarrow \text{prod-pmf } \{..<s\} (\lambda \cdot. \text{pmf-of-set } (\text{bounded-degree-polynomials } (\text{mod-ring } p) \ 2))$ ;  
   return-pmf  $(s, t, p, r, h, (\lambda \cdot \in \{0..<s\}. \{\}))$   
 }

**fun**  $f0\text{-update} :: \text{nat} \Rightarrow f0\text{-state} \Rightarrow f0\text{-state pmf}$  **where**  
 $f0\text{-update } x (s, t, p, r, h, \text{sketch}) =$   
 return-pmf  $(s, t, p, r, h, \lambda i \in \{..<s\}. \text{least } t (\text{insert } (\text{float-of } (\text{truncate-down } r (\text{hash } p \ x \ (h \ i)))) (\text{sketch } i)))$

**fun**  $f0\text{-result} :: f0\text{-state} \Rightarrow \text{rat pmf}$  **where**  
 $f0\text{-result } (s, t, p, r, h, \text{sketch}) = \text{return-pmf } (\text{median } s (\lambda i \in \{..<s\}. \text{if card } (\text{sketch } i) < t \text{ then of-nat } (\text{card } (\text{sketch } i)) \text{ else } \text{rat-of-nat } t * \text{rat-of-nat } p / \text{rat-of-float } (\text{Max } (\text{sketch } i))))$   
 $))$

**fun**  $f0\text{-space-usage} :: (\text{nat} \times \text{rat} \times \text{rat}) \Rightarrow \text{real}$  **where**  
 $f0\text{-space-usage } (n, \varepsilon, \delta) =$   
 let  $s = \text{nat } \lceil -18 * \ln (\text{real-of-rat } \varepsilon) \rceil$  in  
 let  $r = \text{nat } (4 * \lceil \log 2 (1 / \text{real-of-rat } \delta) \rceil + 23)$  in  
 let  $t = \text{nat } \lceil 80 / (\text{real-of-rat } \delta)^2 \rceil$  in  
 $6 +$   
 $2 * \log 2 (\text{real } s + 1) +$   
 $2 * \log 2 (\text{real } t + 1) +$   
 $2 * \log 2 (\text{real } n + 21) +$   
 $2 * \log 2 (\text{real } r + 1) +$

$real\ s * (5 + 2 * \log 2\ (21 + real\ n) +$   
 $real\ t * (13 + 4 * r + 2 * \log 2\ (\log 2\ (real\ n + 13))))$

**definition** *encode-f0-state* :: *f0-state*  $\Rightarrow$  *bool list option* **where**

*encode-f0-state* =  
 $N_e \bowtie_e (\lambda s.$   
 $N_e \times_e ($   
 $N_e \bowtie_e (\lambda p.$   
 $N_e \times_e ($   
 $([0..<s] \rightarrow_e (P_e\ p\ 2)) \times_e$   
 $([0..<s] \rightarrow_e (S_e\ F_e))))$

**lemma** *inj-on encode-f0-state* (*dom encode-f0-state*)  
 $\langle proof \rangle$

**context**

**fixes**  $\varepsilon\ \delta :: rat$   
**fixes**  $n :: nat$   
**fixes**  $as :: nat\ list$   
**fixes**  $result$   
**assumes**  $\varepsilon\text{-range}: \varepsilon \in \{0 < .. < 1\}$   
**assumes**  $\delta\text{-range}: \delta \in \{0 < .. < 1\}$   
**assumes**  $as\text{-range}: set\ as \subseteq \{.. < n\}$   
**defines**  $result \equiv fold\ (\lambda a\ state. state \gg= f0\text{-update}\ a)\ as\ (f0\text{-init}\ \delta\ \varepsilon\ n) \gg=$   
 $f0\text{-result}$   
**begin**

**private definition** *t* **where**  $t = nat\ \lceil 80 / (real\text{-of-rat}\ \delta)^2 \rceil$

**private lemma** *t-gt-0*:  $t > 0$   $\langle proof \rangle$  **definition** *s* **where**  $s = nat\ \lceil -(18 * \ln$   
 $(real\text{-of-rat}\ \varepsilon)) \rceil$

**private lemma** *s-gt-0*:  $s > 0$   $\langle proof \rangle$  **definition** *p* **where**  $p = prime\text{-above}\ (max$   
 $n\ 19)$

**private lemma** *p-prime:Factorial-Ring.prime* *p*

$\langle proof \rangle$  **lemma** *p-ge-18*:  $p \geq 18$   
 $\langle proof \rangle$  **lemma** *p-gt-0*:  $p > 0$   $\langle proof \rangle$  **lemma** *p-gt-1*:  $p > 1$   $\langle proof \rangle$  **lemma** *n-le-p*:  
 $n \leq p$   
 $\langle proof \rangle$  **lemma** *p-le-n*:  $p \leq 2 * n + 40$   
 $\langle proof \rangle$  **lemma** *as-lt-p*:  $\bigwedge x. x \in set\ as \implies x < p$   
 $\langle proof \rangle$  **lemma** *as-subset-p*:  $set\ as \subseteq \{.. < p\}$   
 $\langle proof \rangle$  **definition** *r* **where**  $r = nat\ (4 * \lceil \log 2\ (1 / real\text{-of-rat}\ \delta) \rceil + 23)$

**private lemma** *r-bound*:  $4 * \log 2\ (1 / real\text{-of-rat}\ \delta) + 23 \leq r$

$\langle proof \rangle$  **lemma** *r-ge-23*:  $r \geq 23$

$\langle proof \rangle$  **lemma** *two-pow-r-le-1*:  $0 < 1 - 2\ powr - real\ r$

$\langle proof \rangle$

**interpretation** *carter-wegman-hash-family mod-ring* *p* 2

**rewrites** *ring.hash* (*mod-ring* *p*) = *Frequency-Moment-0.hash* *p*

$\langle \text{proof} \rangle$  **definition** *tr-hash* **where**  $\text{tr-hash } x \ \omega = \text{truncate-down } r \ (\text{hash } x \ \omega)$

**private definition** *sketch-rv* **where**  
 $\text{sketch-rv } \omega = \text{least } t \ ((\lambda x. \text{float-of } (\text{tr-hash } x \ \omega)) \text{ ' set as})$

**private definition** *estimate*  
**where**  $\text{estimate } S = (\text{if card } S < t \text{ then of-nat } (\text{card } S) \text{ else of-nat } t * \text{of-nat } p / \text{rat-of-float } (\text{Max } S))$

**private definition** *sketch-rv'* **where**  $\text{sketch-rv}' \ \omega = \text{least } t \ ((\lambda x. \text{tr-hash } x \ \omega) \text{ ' set as})$

**private definition** *estimate'* **where**  $\text{estimate}' \ S = (\text{if card } S < t \text{ then real } (\text{card } S) \text{ else real } t * \text{real } p / \text{Max } S)$

**private definition**  $\Omega_0$  **where**  $\Omega_0 = \text{prod-pmf } \{..<s\} \ (\lambda-. \text{pmf-of-set space})$

**private lemma** *f0-alg-sketch*:  
**defines**  $\text{sketch} \equiv \text{fold } (\lambda a \text{ state. state } \gg= \text{f0-update } a) \text{ as } (\text{f0-init } \delta \ \varepsilon \ n)$   
**shows**  $\text{sketch} = \text{map-pmf } (\lambda x. (s, t, p, r, x, \lambda i \in \{..<s\}. \text{sketch-rv } (x \ i))) \ \Omega_0$   
 $\langle \text{proof} \rangle$  **lemma** *card-nat-in-ball*:  
**fixes**  $x :: \text{nat}$   
**fixes**  $q :: \text{real}$   
**assumes**  $q \geq 0$   
**defines**  $A \equiv \{k. \text{abs } (\text{real } x - \text{real } k) \leq q \wedge k \neq x\}$   
**shows**  $\text{real } (\text{card } A) \leq 2 * q \text{ and finite } A$   
 $\langle \text{proof} \rangle$  **lemma** *prob-degree-lt-1*:  
 $\text{prob } \{\omega. \text{degree } \omega < 1\} \leq 1 / \text{real } p$   
 $\langle \text{proof} \rangle$  **lemma** *collision-prob*:  
**assumes**  $c \geq 1$   
**shows**  $\text{prob } \{\omega. \exists x \in \text{set as. } \exists y \in \text{set as. } x \neq y \wedge \text{tr-hash } x \ \omega \leq c \wedge \text{tr-hash } x \ \omega = \text{tr-hash } y \ \omega\} \leq$   
 $(5/2) * (\text{real } (\text{card } (\text{set as})))^2 * c^2 * 2 \text{ powr } -(\text{real } r) / (\text{real } p)^2 + 1 / \text{real } p$   
**(is prob**  $\{\omega. ?l \ \omega\} \leq ?r1 + ?r2)$   
 $\langle \text{proof} \rangle$  **lemma** *of-bool-square*:  $(\text{of-bool } x)^2 = ((\text{of-bool } x)::\text{real})$   
 $\langle \text{proof} \rangle$  **definition** *Q* **where**  $Q \ y \ \omega = \text{card } \{x \in \text{set as. int } (\text{hash } x \ \omega) < y\}$

**private definition** *m* **where**  $m = \text{card } (\text{set as})$

**private lemma**  
**assumes**  $a \geq 0$   
**assumes**  $a \leq \text{int } p$   
**shows**  $\text{exp-Q: expectation } (\lambda \omega. \text{real } (Q \ a \ \omega)) = \text{real } m * (\text{of-int } a) / p$   
**and**  $\text{var-Q: variance } (\lambda \omega. \text{real } (Q \ a \ \omega)) \leq \text{real } m * (\text{of-int } a) / p$   
 $\langle \text{proof} \rangle$  **lemma** *t-bound*:  $t \leq 81 / (\text{real-of-rat } \delta)^2$   
 $\langle \text{proof} \rangle$  **lemma** *t-r-bound*:  
 $18 * 40 * (\text{real } t)^2 * 2 \text{ powr } (-\text{real } r) \leq 1$   
 $\langle \text{proof} \rangle$  **lemma** *m-eq-F-0*:  $\text{real } m = \text{of-rat } (F \ 0 \ \text{as})$   
 $\langle \text{proof} \rangle$  **lemma** *estimate'-bounds*:  
 $\text{prob } \{\omega. \text{of-rat } \delta * \text{real-of-rat } (F \ 0 \ \text{as}) < |\text{estimate}' (\text{sketch-rv}' \ \omega) - \text{of-rat } (F \ 0$

$as)|\} \leq 1/3$   
 $\langle proof \rangle$  **lemma** *median-bounds*:  
 $\mathcal{P}(\omega \text{ in measure-pmf } \Omega_0. |\text{median } s (\lambda i. \text{estimate } (\text{sketch-rv } (\omega \ i))) - F \ 0 \ as| \leq$   
 $\delta * F \ 0 \ as) \geq 1 - \text{real-of-rat } \varepsilon$   
 $\langle proof \rangle$

**lemma** *f0-alg-correct'*:  
 $\mathcal{P}(\omega \text{ in measure-pmf result. } |\omega - F \ 0 \ as| \leq \delta * F \ 0 \ as) \geq 1 - \text{of-rat } \varepsilon$   
 $\langle proof \rangle$  **lemma** *f-subset*:  
**assumes**  $g \text{ ' } A \subseteq h \text{ ' } B$   
**shows**  $(\lambda x. f \ (g \ x)) \text{ ' } A \subseteq (\lambda x. f \ (h \ x)) \text{ ' } B$   
 $\langle proof \rangle$

**lemma** *f0-exact-space-usage'*:  
**defines**  $\Omega \equiv \text{fold } (\lambda a \text{ state. state } \gg f0\text{-update } a) \text{ as } (f0\text{-init } \delta \ \varepsilon \ n)$   
**shows**  $AE \ \omega \text{ in } \Omega. \text{bit-count } (\text{encode-f0-state } \omega) \leq f0\text{-space-usage } (n, \varepsilon, \delta)$   
 $\langle proof \rangle$

**end**

Main results of this section:

**theorem** *f0-alg-correct*:  
**assumes**  $\varepsilon \in \{0 < .. < 1\}$   
**assumes**  $\delta \in \{0 < .. < 1\}$   
**assumes**  $\text{set } as \subseteq \{.. < n\}$   
**defines**  $\Omega \equiv \text{fold } (\lambda a \text{ state. state } \gg f0\text{-update } a) \text{ as } (f0\text{-init } \delta \ \varepsilon \ n) \gg f0\text{-result}$   
**shows**  $\mathcal{P}(\omega \text{ in measure-pmf } \Omega. |\omega - F \ 0 \ as| \leq \delta * F \ 0 \ as) \geq 1 - \text{of-rat } \varepsilon$   
 $\langle proof \rangle$

**theorem** *f0-exact-space-usage*:  
**assumes**  $\varepsilon \in \{0 < .. < 1\}$   
**assumes**  $\delta \in \{0 < .. < 1\}$   
**assumes**  $\text{set } as \subseteq \{.. < n\}$   
**defines**  $\Omega \equiv \text{fold } (\lambda a \text{ state. state } \gg f0\text{-update } a) \text{ as } (f0\text{-init } \delta \ \varepsilon \ n)$   
**shows**  $AE \ \omega \text{ in } \Omega. \text{bit-count } (\text{encode-f0-state } \omega) \leq f0\text{-space-usage } (n, \varepsilon, \delta)$   
 $\langle proof \rangle$

**theorem** *f0-asymptotic-space-complexity*:  
 $f0\text{-space-usage} \in O[\text{at-top } \times_F \text{at-right } 0 \times_F \text{at-right } 0](\lambda(n, \varepsilon, \delta). \ln(1 / \text{of-rat } \varepsilon) * \\ (\ln(\text{real } n) + 1 / (\text{of-rat } \delta)^2 * (\ln(\ln(\text{real } n)) + \ln(1 / \text{of-rat } \delta))))$   
**(is -**  $\in O[?F](?rhs)$ **)**  
 $\langle proof \rangle$

**end**

## 8 Frequency Moment 2

**theory** *Frequency-Moment-2*



```

imports
  Universal-Hash-Families.Carter-Wegman-Hash-Family
  Universal-Hash-Families.Universal-Hash-Families-More-Finite-Fields
  Equivalence-Relation-Enumeration.Equivalence-Relation-Enumeration
  Landau-Ext
  Median-Method.Median
  Probability-Ext
  Product-PMF-Ext
  Frequency-Moments
begin

hide-const (open) Discrete-Topology.discrete
hide-const (open) Isolated.discrete

This section contains a formalization of the algorithm for the second frequency moment. It is based on the algorithm described in [1, §2.2]. The only difference is that the algorithm is adapted to work with prime field of odd order, which greatly reduces the implementation complexity.

fun f2-hash where
  f2-hash p h k = (if even (ring.hash (mod-ring p) k h) then int p - 1 else - int p - 1)

type-synonym f2-state = nat × nat × nat × (nat × nat ⇒ nat list) × (nat × nat ⇒ int)

fun f2-init :: rat ⇒ rat ⇒ nat ⇒ f2-state pmf where
  f2-init δ ε n =
    do {
      let s1 = nat ⌈6 / δ2⌉;
      let s2 = nat ⌈-(18 * ln (real-of-rat ε))⌉;
      let p = prime-above (max n 3);
      h ← prod-pmf ({..s1} × {..s2}) (λ-. pmf-of-set (bounded-degree-polynomials (mod-ring p) 4));
      return-pmf (s1, s2, p, h, (λ- ∈ {..s1} × {..s2}. (0 :: int)))
    }

fun f2-update :: nat ⇒ f2-state ⇒ f2-state pmf where
  f2-update x (s1, s2, p, h, sketch) =
    return-pmf (s1, s2, p, h, λi ∈ {..s1} × {..s2}. f2-hash p (h i) x + sketch i)

fun f2-result :: f2-state ⇒ rat pmf where
  f2-result (s1, s2, p, h, sketch) =
    return-pmf (median s2 (λi2 ∈ {..s2}.
      (∑ i1 ∈ {..s1} . (rat-of-int (sketch (i1, i2)))2) / (((rat-of-nat p)2 - 1) *
      rat-of-nat s1)))

fun f2-space-usage :: (nat × nat × rat × rat) ⇒ real where
  f2-space-usage (n, m, ε, δ) = (
    let s1 = nat ⌈6 / δ2⌉ in

```

$let\ s_2 = nat\ \lceil -(18 * ln\ (real-of-rat\ \varepsilon)) \rceil\ in$   
 $3 +$   
 $2 * log\ 2\ (s_1 + 1) +$   
 $2 * log\ 2\ (s_2 + 1) +$   
 $2 * log\ 2\ (9 + 2 * real\ n) +$   
 $s_1 * s_2 * (5 + 4 * log\ 2\ (8 + 2 * real\ n) + 2 * log\ 2\ (real\ m * (18 + 4 * real$   
 $n) + 1\ )))$

**definition** *encode-f2-state* :: *f2-state*  $\Rightarrow$  *bool list option* **where**

*encode-f2-state* =  
 $N_e \bowtie_e (\lambda s_1.$   
 $N_e \bowtie_e (\lambda s_2.$   
 $N_e \bowtie_e (\lambda p.$   
 $(List.product\ [0..<s_1]\ [0..<s_2]\ \rightarrow_e\ P_e\ p\ 4)\ \times_e$   
 $(List.product\ [0..<s_1]\ [0..<s_2]\ \rightarrow_e\ I_e))))$

**lemma** *inj-on encode-f2-state* (*dom encode-f2-state*)  
 $\langle proof \rangle$

**context**

**fixes**  $\varepsilon\ \delta :: rat$   
**fixes**  $n :: nat$   
**fixes**  $as :: nat\ list$   
**fixes** *result*  
**assumes**  $\varepsilon\text{-range}: \varepsilon \in \{0 < .. < 1\}$   
**assumes**  $\delta\text{-range}: \delta > 0$   
**assumes**  $as\text{-range}: set\ as \subseteq \{..<n\}$   
**defines**  $result \equiv fold\ (\lambda a\ state. state \ggg f2\text{-update}\ a)\ as\ (f2\text{-init}\ \delta\ \varepsilon\ n) \ggg$   
*f2-result*  
**begin**

**private definition**  $s_1$  **where**  $s_1 = nat\ \lceil 6 / \delta^2 \rceil$

**lemma** *s1-gt-0*:  $s_1 > 0$   
 $\langle proof \rangle$  **definition**  $s_2$  **where**  $s_2 = nat\ \lceil -(18 * ln\ (real-of-rat\ \varepsilon)) \rceil$

**lemma** *s2-gt-0*:  $s_2 > 0$   
 $\langle proof \rangle$  **definition**  $p$  **where**  $p = prime\text{-above}\ (max\ n\ 3)$

**lemma** *p-prime*: *Factorial-Ring.prime*  $p$   
 $\langle proof \rangle$

**lemma** *p-ge-3*:  $p \geq 3$   
 $\langle proof \rangle$

**lemma** *p-gt-0*:  $p > 0$   $\langle proof \rangle$

**lemma** *p-gt-1*:  $p > 1$   $\langle proof \rangle$

**lemma** *p-ge-n*:  $p \geq n$   $\langle$ proof $\rangle$

**interpretation** *carter-wegman-hash-family mod-ring p 4*  
 $\langle$ proof $\rangle$

**definition** *sketch* **where**  $sketch = fold (\lambda a \ state. state \gg= f2\text{-}update\ a) \ as \ (f2\text{-}init\ \delta \ \varepsilon \ n)$

**private definition**  $\Omega$  **where**  $\Omega = prod\text{-}pmf \ (\{..<s_1\} \times \{..<s_2\}) \ (\lambda-. \ pmf\text{-}of\text{-}set\ space)$

**private definition**  $\Omega_p$  **where**  $\Omega_p = measure\text{-}pmf \ \Omega$

**private definition** *sketch-rv* **where**  $sketch\text{-}rv \ \omega = of\text{-}int \ (sum\text{-}list \ (map \ (f2\text{-}hash \ p \ \omega) \ as)) \wedge^2$

**private definition** *mean-rv* **where**  $mean\text{-}rv \ \omega = (\lambda i_2. \ (\sum i_1 = 0..<s_1. \ sketch\text{-}rv \ (\omega \ i_1, \ i_2))) / (((of\text{-}nat \ p)^2 - 1) * of\text{-}nat \ s_1))$

**private definition** *result-rv* **where**  $result\text{-}rv \ \omega = median \ s_2 \ (\lambda i_2 \in \{..<s_2\}. \ mean\text{-}rv \ \omega \ i_2)$

**lemma** *mean-rv-alg-sketch*:

$sketch = \Omega \gg= (\lambda \omega. \ return\text{-}pmf \ (s_1, \ s_2, \ p, \ \omega, \ \lambda i \in \{..<s_1\} \times \{..<s_2\}. \ sum\text{-}list \ (map \ (f2\text{-}hash \ p \ (\omega \ i)) \ as)))$   
 $\langle$ proof $\rangle$

**lemma** *distr*:  $result = map\text{-}pmf \ result\text{-}rv \ \Omega$

$\langle$ proof $\rangle$  **lemma** *f2-hash-pow-exp*:

**assumes**  $k < p$

**shows**

$expectation \ (\lambda \omega. \ real\text{-}of\text{-}int \ (f2\text{-}hash \ p \ \omega \ k) \wedge^m) =$   
 $((real \ p - 1) \wedge^m * (real \ p + 1) + (- \ real \ p - 1) \wedge^m * (real \ p - 1)) / (2 * real \ p)$   
 $\langle$ proof $\rangle$

**lemma**

**shows**  $var\text{-}sketch\text{-}rv\text{:}variance \ sketch\text{-}rv \leq 2 * (real\text{-}of\text{-}rat \ (F \ 2 \ as) \wedge^2) * ((real \ p)^2 - 1)^2$  **(is ?A)**

**and**  $exp\text{-}sketch\text{-}rv\text{:}expectation \ sketch\text{-}rv = real\text{-}of\text{-}rat \ (F \ 2 \ as) * ((real \ p)^2 - 1)$  **(is ?B)**

$\langle$ proof $\rangle$

**lemma** *space-omega-1* [simp]:  $Sigma\text{-}Algebra.space \ \Omega_p = UNIV$

$\langle$ proof $\rangle$

**interpretation**  $\Omega$ : *prob-space*  $\Omega_p$

$\langle$ proof $\rangle$

**lemma** *integrable-Ω*:

**fixes**  $f :: ((nat \times nat) \Rightarrow (nat \ list)) \Rightarrow real$

**shows** *integrable*  $\Omega_p \ f$

$\langle$ proof $\rangle$

**lemma** *sketch-rv-exp*:

**assumes**  $i_2 < s_2$

**assumes**  $i_1 \in \{0..<s_1\}$

**shows**  $\Omega.\text{expectation } (\lambda\omega. \text{sketch-rv } (\omega (i_1, i_2))) = \text{real-of-rat } (F \ 2 \ as) * ((\text{real } p)^2 - 1)$   
 $\langle \text{proof} \rangle$

**lemma** *sketch-rv-var*:

**assumes**  $i_2 < s_2$

**assumes**  $i_1 \in \{0..<s_1\}$

**shows**  $\Omega.\text{variance } (\lambda\omega. \text{sketch-rv } (\omega (i_1, i_2))) \leq 2 * (\text{real-of-rat } (F \ 2 \ as))^2 * ((\text{real } p)^2 - 1)^2$   
 $\langle \text{proof} \rangle$

**lemma** *mean-rv-exp*:

**assumes**  $i < s_2$

**shows**  $\Omega.\text{expectation } (\lambda\omega. \text{mean-rv } \omega \ i) = \text{real-of-rat } (F \ 2 \ as)$   
 $\langle \text{proof} \rangle$

**lemma** *mean-rv-var*:

**assumes**  $i < s_2$

**shows**  $\Omega.\text{variance } (\lambda\omega. \text{mean-rv } \omega \ i) \leq (\text{real-of-rat } (\delta * F \ 2 \ as))^2 / 3$   
 $\langle \text{proof} \rangle$

**lemma** *mean-rv-bounds*:

**assumes**  $i < s_2$

**shows**  $\Omega.\text{prob } \{\omega. \text{real-of-rat } \delta * \text{real-of-rat } (F \ 2 \ as) < |\text{mean-rv } \omega \ i - \text{real-of-rat } (F \ 2 \ as)|\} \leq 1/3$   
 $\langle \text{proof} \rangle$

**lemma** *f2-alg-correct'*:

$\mathcal{P}(\omega \text{ in measure-pmf result. } |\omega - F \ 2 \ as| \leq \delta * F \ 2 \ as) \geq 1 - \text{of-rat } \varepsilon$   
 $\langle \text{proof} \rangle$

**lemma** *f2-exact-space-usage'*:

$AE \ \omega \text{ in sketch. } \text{bit-count } (\text{encode-f2-state } \omega) \leq \text{f2-space-usage } (n, \text{length } as, \varepsilon, \delta)$   
 $\langle \text{proof} \rangle$

**end**

Main results of this section:

**theorem** *f2-alg-correct*:

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

**assumes**  $\delta > 0$

**assumes**  $\text{set } as \subseteq \{..<n\}$

**defines**  $\Omega \equiv \text{fold } (\lambda a \text{ state. state } \gg \text{f2-update } a) \text{ as } (\text{f2-init } \delta \ \varepsilon \ n) \gg \text{f2-result}$

**shows**  $\mathcal{P}(\omega \text{ in measure-pmf } \Omega. |\omega - F \ 2 \ as| \leq \delta * F \ 2 \ as) \geq 1 - \text{of-rat } \varepsilon$

$\langle \text{proof} \rangle$

**theorem** *f2-exact-space-usage*:

**assumes**  $\varepsilon \in \{0 < .. < 1\}$   
**assumes**  $\delta > 0$   
**assumes**  $set\ as \subseteq \{.. < n\}$   
**defines**  $M \equiv fold\ (\lambda a\ state.\ state \gg= f2\text{-}update\ a)\ as\ (f2\text{-}init\ \delta\ \varepsilon\ n)$   
**shows**  $AE\ \omega\ in\ M.\ bit\text{-}count\ (encode\text{-}f2\text{-}state\ \omega) \leq f2\text{-}space\text{-}usage\ (n,\ length\ as,\ \varepsilon,\ \delta)$   
*<proof>*

**theorem** *f2-asymptotic-space-complexity*:

$f2\text{-}space\text{-}usage \in O[at\text{-}top \times_F at\text{-}top \times_F at\text{-}right\ 0 \times_F at\text{-}right\ 0](\lambda\ (n,\ m,\ \varepsilon,\ \delta).\ (ln\ (1 / of\text{-}rat\ \varepsilon)) / (of\text{-}rat\ \delta)^2 * (ln\ (real\ n) + ln\ (real\ m))))$   
**(is -  $\in O[?F](?rhs)$ )**  
*<proof>*

**end**

## 9 Frequency Moment $k$

**theory** *Frequency-Moment-k*

**imports**

*Frequency-Moments*  
*Landau-Ext*  
*Lp.Lp*  
*Median-Method.Median*  
*Probability-Ext*  
*Product-PMF-Ext*

**begin**

This section contains a formalization of the algorithm for the  $k$ -th frequency moment. It is based on the algorithm described in [1, §2.1].

**type-synonym**  $fk\text{-}state = nat \times nat \times nat \times nat \times (nat \times nat \Rightarrow (nat \times nat))$

**fun**  $fk\text{-}init :: nat \Rightarrow rat \Rightarrow rat \Rightarrow nat \Rightarrow fk\text{-}state\ pmf$  **where**

$fk\text{-}init\ k\ \delta\ \varepsilon\ n =$   
*do* {  
 $let\ s_1 = nat\ \lceil 3 * real\ k * n\ powr\ (1 - 1 / real\ k) / (real\text{-}of\text{-}rat\ \delta)^2 \rceil;$   
 $let\ s_2 = nat\ \lceil -18 * ln\ (real\text{-}of\text{-}rat\ \varepsilon) \rceil;$   
 $return\text{-}pmf\ (s_1,\ s_2,\ k,\ 0,\ (\lambda\ - \in \{0..<s_1\} \times \{0..<s_2\}.\ (0,0)))$   
}

**fun**  $fk\text{-}update :: nat \Rightarrow fk\text{-}state \Rightarrow fk\text{-}state\ pmf$  **where**

$fk\text{-}update\ a\ (s_1,\ s_2,\ k,\ m,\ r) =$   
*do* {  
 $coins \leftarrow prod\text{-}pmf\ (\{0..<s_1\} \times \{0..<s_2\})\ (\lambda\ .\ bernoulli\text{-}pmf\ (1 / (real\ m + 1)));$   
 $return\text{-}pmf\ (s_1,\ s_2,\ k,\ m + 1,\ \lambda i \in \{0..<s_1\} \times \{0..<s_2\}.\$   
 $if\ coins\ i\ then$   
 $(a, 0)$   
}

```

    else (
      let (x,l) = r i in (x, l + of-bool (x=a))
    )
  )
}

```

**fun** *fk-result* :: *fk-state*  $\Rightarrow$  *rat pmf* **where**  
*fk-result* (*s*<sub>1</sub>, *s*<sub>2</sub>, *k*, *m*, *r*) =  
 return-pmf (median *s*<sub>2</sub> ( $\lambda i_2 \in \{0..<s_2\}$ .  
 ( $\sum i_1 \in \{0..<s_1\}$ . rat-of-nat (let *t* = snd (*r* (*i*<sub>1</sub>, *i*<sub>2</sub>)) + 1 in *m* \* (*t*  $\wedge$  *k* - (*t* - 1)  $\wedge$  *k*))) / (rat-of-nat *s*<sub>1</sub>))  
 )

**lemma** *bernoulli-pmf-1*: *bernoulli-pmf* 1 = *return-pmf* True  
 <proof>

**fun** *fk-space-usage* :: (*nat*  $\times$  *nat*  $\times$  *nat*  $\times$  *rat*  $\times$  *rat*)  $\Rightarrow$  *real* **where**  
*fk-space-usage* (*k*, *n*, *m*,  $\varepsilon$ ,  $\delta$ ) = (  
 let *s*<sub>1</sub> = nat  $\lceil 3 * \text{real } k * (\text{real } n) \text{ powr } (1 - 1 / \text{real } k) / (\text{real-of-rat } \delta)^2 \rceil$  in  
 let *s*<sub>2</sub> = nat  $\lceil -(18 * \ln (\text{real-of-rat } \varepsilon)) \rceil$  in  
 4 +  
 2 \* log 2 (*s*<sub>1</sub> + 1) +  
 2 \* log 2 (*s*<sub>2</sub> + 1) +  
 2 \* log 2 (*real k* + 1) +  
 2 \* log 2 (*real m* + 1) +  
*s*<sub>1</sub> \* *s*<sub>2</sub> \* (2 + 2 \* log 2 (*real n* + 1) + 2 \* log 2 (*real m* + 1)))

**definition** *encode-fk-state* :: *fk-state*  $\Rightarrow$  *bool list option* **where**  
*encode-fk-state* =  
*N*<sub>*e*</sub>  $\bowtie_e$  ( $\lambda s_1$ .  
*N*<sub>*e*</sub>  $\bowtie_e$  ( $\lambda s_2$ .  
*N*<sub>*e*</sub>  $\times_e$   
*N*<sub>*e*</sub>  $\times_e$   
 (*List.product* [*0*..*s*<sub>1</sub>] [*0*..*s*<sub>2</sub>]  $\rightarrow_e$  (*N*<sub>*e*</sub>  $\times_e$  *N*<sub>*e*</sub>))))

**lemma** *inj-on encode-fk-state* (*dom encode-fk-state*)  
 <proof>

This is an intermediate non-parallel form *fk-update* used only in the correctness proof.

**fun** *fk-update-2* :: '*a*  $\Rightarrow$  (*nat*  $\times$  '*a*  $\times$  *nat*)  $\Rightarrow$  (*nat*  $\times$  '*a*  $\times$  *nat*) pmf **where**  
*fk-update-2* *a* (*m*,*x*,*l*) =  
 do {  
 coin  $\leftarrow$  bernoulli-pmf (1/(*real m* + 1));  
 return-pmf (*m* + 1, if coin then (*a*, 0) else (*x*, *l* + of-bool (*x* = *a*)))  
 }

**definition** *sketch* **where** *sketch as i* = (*as* ! *i*, count-list (*drop* (*i* + 1) *as*) (*as* ! *i*))

**lemma** *fk-update-2-distr*:

**assumes**  $as \neq []$

**shows**  $fold (\lambda x s. s \gg= fk\text{-}update\text{-}2\ x)\ as\ (return\text{-}pmf\ (0,0,0)) =$   
 $pmf\text{-}of\text{-}set\ \{..<length\ as\} \gg= (\lambda k. return\text{-}pmf\ (length\ as,\ sketch\ as\ k))$   
 $\langle proof \rangle$

**context**

**fixes**  $\varepsilon\ \delta :: rat$

**fixes**  $n\ k :: nat$

**fixes**  $as$

**assumes**  $k\text{-ge-}1: k \geq 1$

**assumes**  $\varepsilon\text{-range}: \varepsilon \in \{0 < .. < 1\}$

**assumes**  $\delta\text{-range}: \delta > 0$

**assumes**  $as\text{-range}: set\ as \subseteq \{..<n\}$

**begin**

**definition**  $s_1$  **where**  $s_1 = nat\ \lceil 3 * real\ k * (real\ n)\ powr\ (1 - 1 / real\ k) / (real\text{-}of\text{-}rat\ \delta)^2 \rceil$

**definition**  $s_2$  **where**  $s_2 = nat\ \lceil -(18 * ln\ (real\text{-}of\text{-}rat\ \varepsilon)) \rceil$

**definition**  $M_1 = \{(u, v). v < count\text{-}list\ as\ u\}$

**definition**  $\Omega_1 = measure\text{-}pmf\ (pmf\text{-}of\text{-}set\ M_1)$

**definition**  $M_2 = prod\text{-}pmf\ (\{0..<s_1\} \times \{0..<s_2\})\ (\lambda\cdot. pmf\text{-}of\text{-}set\ M_1)$

**definition**  $\Omega_2 = measure\text{-}pmf\ M_2$

**interpretation**  $prob\text{-}space\ \Omega_1$

$\langle proof \rangle$

**interpretation**  $\Omega_2: prob\text{-}space\ \Omega_2$

$\langle proof \rangle$

**lemma** *split-space*:  $(\sum a \in M_1. f\ (snd\ a)) = (\sum u \in set\ as. (\sum v \in \{0..<count\text{-}list\ as\ u\}. f\ v))$

$\langle proof \rangle$

**lemma**

**assumes**  $as \neq []$

**shows**  $fin\text{-}space: finite\ M_1$

**and**  $non\text{-}empty\text{-}space: M_1 \neq \{\}$

**and**  $card\text{-}space: card\ M_1 = length\ as$

$\langle proof \rangle$

**lemma**

**assumes**  $as \neq []$

**shows**  $integrable\text{-}1: integrable\ \Omega_1\ (f :: - \Rightarrow real)$  **and**

$integrable\text{-}2: integrable\ \Omega_2\ (g :: - \Rightarrow real)$

$\langle proof \rangle$

**lemma** *sketch-distr*:

**assumes**  $as \neq []$   
**shows**  $\text{pmf-of-set } \{..<\text{length } as\} \gg (\lambda k. \text{return-pmf } (\text{sketch } as \ k)) = \text{pmf-of-set } M_1$   
 $\langle \text{proof} \rangle$

**lemma** *fk-update-distr*:

$\text{fold } (\lambda x \ s. s \gg \text{fk-update } x) \ as \ (\text{fk-init } k \ \delta \ \varepsilon \ n) =$   
 $\text{prod-pmf } (\{0..<s_1\} \times \{0..<s_2\}) \ (\lambda-. \text{fold } (\lambda x \ s. s \gg \text{fk-update-2 } x) \ as \ (\text{return-pmf } (0,0,0)))$   
 $\gg (\lambda x. \text{return-pmf } (s_1, s_2, k, \text{length } as, \lambda i \in \{0..<s_1\} \times \{0..<s_2\}. \text{snd } (x \ i)))$   
 $\langle \text{proof} \rangle$

**lemma** *power-diff-sum*:

**fixes**  $a \ b :: 'a :: \{\text{comm-ring-1}, \text{power}\}$   
**assumes**  $k > 0$   
**shows**  $a^{\wedge} k - b^{\wedge} k = (a-b) * (\sum i = 0..<k. a^{\wedge} i * b^{\wedge} (k-1-i))$  (is ?lhs = ?rhs)  
 $\langle \text{proof} \rangle$

**lemma** *power-diff-est*:

**assumes**  $k > 0$   
**assumes**  $(a :: \text{real}) \geq b$   
**assumes**  $b \geq 0$   
**shows**  $a^{\wedge} k - b^{\wedge} k \leq (a-b) * k * a^{\wedge}(k-1)$   
 $\langle \text{proof} \rangle$

Specialization of the Hoelder inequality for sums.

**lemma** *Holder-inequality-sum*:

**assumes**  $p > (0::\text{real}) \ q > 0 \ 1/p + 1/q = 1$   
**assumes** *finite*  $A$   
**shows**  $|\sum x \in A. f \ x * g \ x| \leq (\sum x \in A. |f \ x| \ \text{powr } p) \ \text{powr } (1/p) * (\sum x \in A. |g \ x| \ \text{powr } q) \ \text{powr } (1/q)$   
 $\langle \text{proof} \rangle$

**lemma** *real-count-list-pos*:

**assumes**  $x \in \text{set } as$   
**shows**  $\text{real } (\text{count-list } as \ x) > 0$   
 $\langle \text{proof} \rangle$

**lemma** *fk-estimate*:

**assumes**  $as \neq []$   
**shows**  $\text{length } as * \text{of-rat } (F \ (2*k-1) \ as) \leq n \ \text{powr } (1 - 1 / \text{real } k) * (\text{of-rat } (F \ k \ as))^2$   
 (is ?lhs  $\leq$  ?rhs)  
 $\langle \text{proof} \rangle$

**definition** *result*

**where**  $\text{result } a = \text{of-nat } (\text{length } as) * \text{of-nat } (\text{Suc } (\text{snd } a)^{\wedge} k - \text{snd } a^{\wedge} k)$



**lemma** *result-exp-1*:

**assumes**  $as \neq []$

**shows**  $expectation\ result = real-of-rat\ (F\ k\ as)$

*<proof>*

**lemma** *result-var-1*:

**assumes**  $as \neq []$

**shows**  $variance\ result \leq (of-rat\ (F\ k\ as))^2 * k * n\ powr\ (1 - 1 / real\ k)$

*<proof>*

**theorem** *fk-alg-sketch*:

**assumes**  $as \neq []$

**shows**  $fold\ (\lambda a\ state. state \gg= fk-update\ a)\ as\ (fk-init\ k\ \delta\ \varepsilon\ n) =$   
 $map-pmf\ (\lambda x. (s_1, s_2, k, length\ as, x))\ M_2\ (is\ ?lhs = ?rhs)$

*<proof>*

**definition** *mean-rv*

**where**  $mean-rv\ \omega\ i_2 = (\sum\ i_1 = 0..<s_1. result\ (\omega\ (i_1, i_2))) / of-nat\ s_1$

**definition** *median-rv*

**where**  $median-rv\ \omega = median\ s_2\ (\lambda i_2. mean-rv\ \omega\ i_2)$

**lemma** *fk-alg-correct'*:

**defines**  $M \equiv fold\ (\lambda a\ state. state \gg= fk-update\ a)\ as\ (fk-init\ k\ \delta\ \varepsilon\ n) \gg= fk-result$

**shows**  $\mathcal{P}(\omega\ in\ measure-pmf\ M. |\omega - F\ k\ as| \leq \delta * F\ k\ as) \geq 1 - of-rat\ \varepsilon$

*<proof>*

**lemma** *fk-exact-space-usage'*:

**defines**  $M \equiv fold\ (\lambda a\ state. state \gg= fk-update\ a)\ as\ (fk-init\ k\ \delta\ \varepsilon\ n)$

**shows**  $AE\ \omega\ in\ M. bit-count\ (encode-fk-state\ \omega) \leq fk-space-usage\ (k, n, length\ as, \varepsilon, \delta)$

**(is**  $AE\ \omega\ in\ M. (- \leq ?rhs)$ **)**

*<proof>*

**end**

Main results of this section:

**theorem** *fk-alg-correct*:

**assumes**  $k \geq 1$

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

**assumes**  $\delta > 0$

**assumes**  $set\ as \subseteq \{..<n\}$

**defines**  $M \equiv fold\ (\lambda a\ state. state \gg= fk-update\ a)\ as\ (fk-init\ k\ \delta\ \varepsilon\ n) \gg= fk-result$

**shows**  $\mathcal{P}(\omega\ in\ measure-pmf\ M. |\omega - F\ k\ as| \leq \delta * F\ k\ as) \geq 1 - of-rat\ \varepsilon$

*<proof>*

**theorem** *fk-exact-space-usage*:

**assumes**  $k \geq 1$

```

assumes  $\varepsilon \in \{0 < \dots < 1\}$ 
assumes  $\delta > 0$ 
assumes  $set\ as \subseteq \{.. < n\}$ 
defines  $M \equiv fold\ (\lambda a\ state.\ state \gg= fk\text{-}update\ a)\ as\ (fk\text{-}init\ k\ \delta\ \varepsilon\ n)$ 
shows  $AE\ \omega\ in\ M.\ bit\text{-}count\ (encode\text{-}fk\text{-}state\ \omega) \leq fk\text{-}space\text{-}usage\ (k,\ n,\ length\ as,\ \varepsilon,\ \delta)$ 
 $\langle proof \rangle$ 

```

**theorem** *fk-asymptotic-space-complexity:*

```

   $fk\text{-}space\text{-}usage \in$ 
   $O[at\text{-}top \times_F at\text{-}top \times_F at\text{-}top \times_F at\text{-}right\ (0::rat) \times_F at\text{-}right\ (0::rat)](\lambda\ (k,\ n,$ 
   $m,\ \varepsilon,\ \delta).$ 
   $real\ k * real\ n\ powr\ (1 - 1 / real\ k) / (of\text{-}rat\ \delta)^2 * (\ln\ (1 / of\text{-}rat\ \varepsilon)) * (\ln\ (real$ 
   $n) + \ln\ (real\ m)))$ 
   $(is\ - \in O[?F](?rhs))$ 
 $\langle proof \rangle$ 

```

**end**

## A Informal proof of correctness for the $F_0$ algorithm

This appendix contains a detailed informal proof for the new Rounding-KMV algorithm that approximates  $F_0$  introduced in Section 7 for reference. It follows the same reasoning as the formalized proof.

Because of the amplification result about medians (see for example [1, §2.1]) it is enough to show that each of the estimates the median is taken from is within the desired interval with success probability  $\frac{2}{3}$ . To verify the latter, let  $a_1, \dots, a_m$  be the stream elements, where we assume that the elements are a subset of  $\{0, \dots, n-1\}$  and  $0 < \delta < 1$  be the desired relative accuracy. Let  $p$  be the smallest prime such that  $p \geq \max(n, 19)$  and let  $h$  be a random polynomial over  $GF(p)$  with degree strictly less than 2. The algorithm also introduces the internal parameters  $t, r$  defined by:

$$t := \lceil 80\delta^{-2} \rceil \qquad r := 4 \log_2 \lceil \delta^{-1} \rceil + 23$$

The estimate the algorithm obtains is  $R$ , defined using:

$$H := \{ \lfloor h(a) \rfloor_r \mid a \in A \} \qquad R := \begin{cases} tp(\min_t(H))^{-1} & \text{if } |H| \geq t \\ |H| & \text{otherwise,} \end{cases}$$

where  $A := \{a_1, \dots, a_m\}$ ,  $\min_t(H)$  denotes the  $t$ -th smallest element of  $H$  and  $\lfloor x \rfloor_r$  denotes the largest binary floating point number smaller or equal to  $x$  with a mantissa that requires at most  $r$  bits to represent.<sup>1</sup> With these

---

<sup>1</sup>This rounding operation is called *truncate-down* in Isabelle, it is defined in `HOL-Library.Float`.

definitions, it is possible to state the main theorem as:

$$P(|R - F_0| \leq \delta |F_0|) \geq \frac{2}{3}.$$

which is shown separately in the following two subsections for the cases  $F_0 \geq t$  and  $F_0 < t$ .

### A.1 Case $F_0 \geq t$

Let us introduce:

$$H^* := \{h(a) | a \in A\}^\# \quad R^* := tp \left( \min_t^\#(H^*) \right)^{-1}$$

These definitions are modified versions of the definitions for  $H$  and  $R$ : The set  $H^*$  is a multiset, this means that each element also has a multiplicity, counting the number of *distinct* elements of  $A$  being mapped by  $h$  to the same value. Note that by definition:  $|H^*| = |A|$ . Similarly the operation  $\min_t^\#$  obtains the  $t$ -th element of the multiset  $H$  (taking multiplicities into account). Note also that there is no rounding operation  $\lfloor \cdot \rfloor_r$  in the definition of  $H^*$ . The key reason for the introduction of these alternative versions of  $H, R$  is that it is easier to show probabilistic bounds on the distances  $|R^* - F_0|$  and  $|R^* - R|$  as opposed to  $|R - F_0|$  directly. In particular the plan is to show:

$$P(|R^* - F_0| > \delta' F_0) \leq \frac{2}{9}, \text{ and} \quad (1)$$

$$P\left(|R^* - F_0| \leq \delta' F_0 \wedge |R - R^*| > \frac{\delta}{4} F_0\right) \leq \frac{1}{9} \quad (2)$$

where  $\delta' := \frac{3}{4}\delta$ . I.e. the probability that  $R^*$  has not the relative accuracy of  $\frac{3}{4}\delta$  is less than  $\frac{2}{9}$  and the probability that assuming  $R^*$  has the relative accuracy of  $\frac{3}{4}\delta$  but that  $R$  deviates by more than  $\frac{1}{4}\delta F_0$  is at most  $\frac{1}{9}$ . Hence, the probability that neither of these events happen is at least  $\frac{2}{3}$  but in that case:

$$|R - F_0| \leq |R - R^*| + |R^* - F_0| \leq \frac{\delta}{4} F_0 + \frac{3\delta}{4} F_0 = \delta F_0. \quad (3)$$

Thus we only need to show [Equation 1](#) and [2](#). For the verification of [Equation 1](#) let

$$Q(u) = |\{h(a) < u \mid a \in A\}|$$

and observe that  $\min_t^\#(H^*) < u$  if  $Q(u) \geq t$  and  $\min_t^\#(H^*) \geq v$  if  $Q(v) \leq t - 1$ . To see why this is true note that, if at least  $t$  elements of  $A$  are mapped by  $h$  below a certain value, then the  $t$ -smallest element must also be within them, and thus also be below that value. And that the opposite direction of this conclusion is also true. Note that this relies on the fact

that  $H^*$  is a multiset and that multiplicities are being taken into account, when computing the  $t$ -th smallest element. Alternatively, it is also possible to write  $Q(u) = \sum_{a \in A} 1_{\{h(a) < u\}}$ <sup>2</sup>, i.e.,  $Q$  is a sum of pairwise independent  $\{0, 1\}$ -valued random variables, with expectation  $\frac{u}{p}$  and variance  $\frac{u}{p} - \frac{u^2}{p^2}$ .<sup>3</sup> Using linearity of expectation and Bienaymé's identity, it follows that  $\text{Var } Q(u) \leq \text{E } Q(u) = |A|up^{-1} = F_0up^{-1}$  for  $u \in \{0, \dots, p\}$ .

For  $v = \left\lfloor \frac{tp}{(1-\delta')F_0} \right\rfloor$  it is possible to conclude:

$$t - 1 \leq \frac{t}{(1-\delta')} - 3\sqrt{\frac{t}{(1-\delta')}} - 1 \leq \frac{F_0v}{p} - 3\sqrt{\frac{F_0v}{p}} \leq \text{E}Q(v) - 3\sqrt{\text{Var}Q(v)}$$

and thus using Tchebyshev's inequality:

$$\begin{aligned} P(R^* < (1-\delta')F_0) &= P\left(\text{rank}_t^\#(H^*) > \frac{tp}{(1-\delta')F_0}\right) \\ &\leq P(\text{rank}_t^\#(H^*) \geq v) = P(Q(v) \leq t-1) \\ &\leq P\left(Q(v) \leq \text{E}Q(v) - 3\sqrt{\text{Var}Q(v)}\right) \leq \frac{1}{9}. \end{aligned} \quad (4)$$

Similarly for  $u = \left\lceil \frac{tp}{(1+\delta')F_0} \right\rceil$  it is possible to conclude:

$$t \geq \frac{t}{(1+\delta')} + 3\sqrt{\frac{t}{(1+\delta')}} + 1 + 1 \geq \frac{F_0u}{p} + 3\sqrt{\frac{F_0u}{p}} \geq \text{E}Q(u) + 3\sqrt{\text{Var}Q(u)}$$

and thus using Tchebyshev's inequality:

$$\begin{aligned} P(R^* > (1+\delta')F_0) &= P\left(\text{rank}_t^\#(H^*) < \frac{tp}{(1+\delta')F_0}\right) \\ &\leq P(\text{rank}_t^\#(H^*) < u) = P(Q(u) \geq t) \\ &\leq P\left(Q(u) \geq \text{E}Q(u) + 3\sqrt{\text{Var}Q(u)}\right) \leq \frac{1}{9}. \end{aligned} \quad (5)$$

Note that Equation 4 and 5 confirm Equation 1. To verify Equation 2, note that

$$\min_t(H) = \lfloor \min_t^\#(H^*) \rfloor_r \quad (6)$$

if there are no collisions, induced by the application of  $\lfloor h(\cdot) \rfloor_r$  on the elements of  $A$ . Even more carefully, note that the equation would remain true,

<sup>2</sup>The notation  $1_A$  is shorthand for the indicator function of  $A$ , i.e.,  $1_A(x) = 1$  if  $x \in A$  and 0 otherwise.

<sup>3</sup>A consequence of  $h$  being chosen uniformly from a 2-independent hash family.

<sup>4</sup>The verification of this inequality is a lengthy but straightforward calculation using the definition of  $\delta'$  and  $t$ .

as long as there are no collision within the smallest  $t$  elements of  $H^*$ . Because Equation 2 needs to be shown only in the case where  $R^* \geq (1 - \delta')F_0$ , i.e., when  $\min_t^\#(H^*) \leq v$ , it is enough to bound the probability of a collision in the range  $[0; v]$ . Moreover Equation 6 implies  $|\min_t(H) - \min_t^\#(H^*)| \leq \max(\min_t^\#(H^*), \min_t(H))2^{-r}$  from which it is possible to derive  $|R^* - R| \leq \frac{\delta}{4}F_0$ . Another important fact is that  $h$  is injective with probability  $1 - \frac{1}{p}$ , this is because  $h$  is chosen uniformly from the polynomials of degree less than 2. If it is a degree 1 polynomial it is a linear function on  $GF(p)$  and thus injective. Because  $p \geq 18$  the probability that  $h$  is not injective can be bounded by  $1/18$ . With these in mind, we can conclude:

$$\begin{aligned}
& P\left(|R^* - F_0| \leq \delta'F_0 \wedge |R - R^*| > \frac{\delta}{4}F_0\right) \\
& \leq P\left(R^* \geq (1 - \delta')F_0 \wedge \min_t^\#(H^*) \neq \min_t(H) \wedge h \text{ inj.}\right) + P(\neg h \text{ inj.}) \\
& \leq P(\exists a \neq b \in A. \lfloor h(a) \rfloor_r = \lfloor h(b) \rfloor_r \leq v \wedge h(a) \neq h(b)) + \frac{1}{18} \\
& \leq \frac{1}{18} + \sum_{a \neq b \in A} P(\lfloor h(a) \rfloor_r = \lfloor h(b) \rfloor_r \leq v \wedge h(a) \neq h(b)) \\
& \leq \frac{1}{18} + \sum_{a \neq b \in A} P(|h(a) - h(b)| \leq v2^{-r} \wedge h(a) \leq v(1 + 2^{-r}) \wedge h(a) \neq h(b)) \\
& \leq \frac{1}{18} + \sum_{a \neq b \in A} \sum_{\substack{a', b' \in \{0, \dots, p-1\} \wedge a' \neq b' \\ |a' - b'| \leq v2^{-r} \wedge a' \leq v(1 + 2^{-r})}} P(h(a) = a')P(h(b) = b') \\
& \leq \frac{1}{18} + \frac{5F_0^2 v^2}{2p^2} 2^{-r} \leq \frac{1}{9}.
\end{aligned}$$

which shows that Equation 2 is true.

## A.2 Case $F_0 < t$

Note that in this case  $|H| \leq F_0 < t$  and thus  $R = |H|$ , hence the goal is to show that:  $P(|H| \neq F_0) \leq \frac{1}{3}$ . The latter can only happen, if there is a collision induced by the application of  $\lfloor h(\cdot) \rfloor_r$ . As before  $h$  is not injective

with probability at most  $\frac{1}{18}$ , hence:

$$\begin{aligned}
& P(|R - F_0| > \delta F_0) \leq P(R \neq F_0) \\
& \leq \frac{1}{18} + P(R \neq F_0 \wedge h \text{ inj.}) \\
& \leq \frac{1}{18} + P(\exists a \neq b \in A. \lfloor h(a) \rfloor_r = \lfloor h(b) \rfloor_r \wedge h \text{ inj.}) \\
& \leq \frac{1}{18} + \sum_{a \neq b \in A} P(\lfloor h(a) \rfloor_r = \lfloor h(b) \rfloor_r \wedge h(a) \neq h(b)) \\
& \leq \frac{1}{18} + \sum_{a \neq b \in A} P(|h(a) - h(b)| \leq p2^{-r} \wedge h(a) \neq h(b)) \\
& \leq \frac{1}{18} + \sum_{a \neq b \in A} \sum_{\substack{a', b' \in \{0, \dots, p-1\} \\ a' \neq b' \wedge |a' - b'| \leq p2^{-r}}} P(h(a) = a')P(h(b) = b') \\
& \leq \frac{1}{18} + F_0^2 2^{-r+1} \leq \frac{1}{18} + t^2 2^{-r+1} \leq \frac{1}{9}.
\end{aligned}$$

Which concludes the proof.  $\square$

## References

- [1] N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences*, 58(1):137–147, 1999.
- [2] Z. Bar-Yossef, T. S. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan. Counting distinct elements in a data stream. In J. D. P. Rolim and S. Vadhan, editors, *Randomization and Approximation Techniques in Computer Science*, pages 1–10. Springer Berlin Heidelberg, 2002.