

Formalization of Randomized Approximation Algorithms for Frequency Moments

Emin Karayel

March 29, 2023

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

1	Preliminary Results	2
2	Frequency Moments	5
3	Ranks, k smallest element and elements	6
4	Landau Symbols	9
5	Probability Spaces	11
6	Indexed Products of Probability Mass Functions	14
7	Frequency Moment 0	16
8	Frequency Moment 2	20
9	Frequency Moment k	25

A Informal proof of correctness for the F_0 algorithm	30
A.1 Case $F_0 \geq t$	31
A.2 Case $F_0 < t$	33

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

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 *suffix* (*drop* x a) (*drop* y a)
 ⟨*proof*⟩

lemma *count-list-card*: *count-list* xs x = *card* { k . k < *length* xs ∧ xs ! k = x }
 ⟨*proof*⟩

lemma *card-gr-1-iff*:
assumes *finite* S x ∈ S y ∈ S x ≠ y
shows *card* S > 1
 ⟨*proof*⟩

lemma *count-list-ge-2-iff*:
assumes y < z
assumes z < *length* xs
assumes xs ! y = xs ! z
shows *count-list* xs (xs ! y) > 1
 ⟨*proof*⟩

Results about multisets and sorting

This is an induction scheme over the distinct elements of a multiset: We can represent each multiset as a sum like: *replicate-mset* n_1 x_1 + *replicate-mset* n_2 x_2 + ... + *replicate-mset* n_k x_k where the x_i are distinct.

lemma *disj-induct-mset*:
assumes P { $\#$ }
assumes $\bigwedge n$ M x . P M \implies $\neg(x \in \# M) \implies n > 0 \implies P$ (M + *replicate-mset* n x)
shows P M
 ⟨*proof*⟩

lemma *prod-mset-conv*:
fixes f :: 'a \Rightarrow 'b::{*comm-monoid-mult*}
shows *prod-mset* (*image-mset* f A) = *prod* (λx . f x \frown (*count* A x)) (*set-mset* A)
 ⟨*proof*⟩

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::{*ring,semiring-1*}
shows *sum-list* (*map* f xs) = ($\sum x \in$ *set* xs . *of-nat* (*count-list* xs x) * f x)
 ⟨*proof*⟩

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

lemma *sorted-sorted-list-of-multiset*: *sorted* (*sorted-list-of-multiset* M)
 ⟨*proof*⟩

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

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

lemma *list-eq-iff*:
assumes $\text{mset } xs = \text{mset } ys$
assumes *sorted* xs
assumes *sorted* ys
shows $xs = ys$
<proof>

lemma *sorted-list-of-multiset-image-commute*:
assumes *mono* f
shows $\text{sorted-list-of-multiset } (\text{image-mset } f M) = \text{map } f (\text{sorted-list-of-multiset } M)$
<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 $\text{truncate-down } r x = \text{truncate-down } r y$
shows $\text{abs } (x - y) \leq \text{max } (\text{abs } x) (\text{abs } y) * 2 \text{ powr } (-\text{real } r)$
<proof>

definition *rat-of-float* :: $\text{float} \Rightarrow \text{rat}$ **where**
 $\text{rat-of-float } f = \text{of-int } (\text{mantissa } f) * (\text{if } \text{exponent } f \geq 0 \text{ then } 2 ^ (\text{nat } (\text{exponent } f)) \text{ else } 1 / 2 ^ (\text{nat } (-\text{exponent } f)))$

lemma *real-of-rat-of-float*: $\text{real-of-rat } (\text{rat-of-float } x) = \text{real-of-float } x$
<proof>

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

<proof>

lemma *truncate-mantissa-bound*:

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

lemma *truncate-float-bit-count*:

bit-count ($F_e(\text{float-of}(\text{truncate-down } r \ x)) \leq 10 + 4 * \text{real } r + 2 * \log 2 (2 + \lfloor \log 2 |x| \rfloor)$)
(**is** *?lhs* \leq *?rhs*)
<proof>

definition *prime-above* :: *nat* \Rightarrow *nat*

where *prime-above* $n = (\text{SOME } x. x \in \{n..(2*n+2)\} \wedge \text{prime } x)$

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 $\exists x \in A. P \ x$
assumes $A \subseteq B$
shows $\exists x \in B. P \ x$
<proof>

lemma

shows *prime-above-prime*: *prime* (*prime-above* n)
and *prime-above-range*: *prime-above* $n \in \{n..(2*n+2)\}$
<proof>

lemma *prime-above-min*: *prime-above* $n \geq 2$

<proof>

lemma *prime-above-lower-bound*: *prime-above* $n \geq n$

<proof>

lemma *prime-above-upper-bound*: *prime-above* $n \leq 2*n+2$

<proof>

end

2 Frequency Moments

theory *Frequency-Moments*

imports

Frequency-Moments-Preliminary-Results

Universal-Hash-Families.Field
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**

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

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

<proof>

lemma *F-gr-0*:

assumes $as \neq []$

shows $F\ k\ as > 0$

<proof>

definition $P_e :: nat \Rightarrow nat \Rightarrow nat\ list \Rightarrow bool\ list\ option$ **where**

$P_e\ p\ n\ f = (if\ p > 1 \wedge f \in\ bounded-degree-polynomials\ (Field.mod-ring\ p)\ n\ then$
 $([0..<n] \rightarrow_e\ Nb_e\ p)\ (\lambda i \in\ \{..<n\}.\ ring.coeff\ (Field.mod-ring\ p)\ f\ i)\ else\ None)$

lemma *poly-encoding*:

is-encoding $(P_e\ p\ n)$

<proof>

lemma *bounded-degree-polynomial-bit-count*:

assumes $p > 1$

assumes $x \in\ bounded-degree-polynomials\ (Field.mod-ring\ p)\ n$

shows $bit-count\ (P_e\ p\ n\ x) \leq ereal\ (real\ n * (\log\ 2\ p + 1))$

<proof>

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 $x \leq y \implies \text{rank-of } x \ S \leq \text{rank-of } y \ S$
 ⟨*proof*⟩

lemma *rank-mono-2*:
assumes *finite S*
shows $S' \subseteq S \implies \text{rank-of } x \ S' \leq \text{rank-of } x \ S$
 ⟨*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** $\text{count-le } x M = \text{size } \{\#y \in \# M. y \leq x\# \}$

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

definition *nth-mset* $:: \text{nat} \Rightarrow ('a :: \text{linorder}) \text{multiset} \Rightarrow 'a$ **where**

$\text{nth-mset } k M = \text{sorted-list-of-multiset } M ! k$

lemma *nth-mset-bound-left:*

assumes $k < \text{size } M$

assumes $\text{count-less } x M \leq k$

shows $x \leq \text{nth-mset } k M$

<proof>

lemma *nth-mset-bound-left-excl:*

assumes $k < \text{size } M$

assumes $\text{count-le } x M \leq k$

shows $x < \text{nth-mset } k M$

<proof>

lemma *nth-mset-bound-right:*

assumes $k < \text{size } M$

assumes $\text{count-le } x M > k$

shows $\text{nth-mset } k M \leq x$

<proof>

lemma *nth-mset-commute-mono:*

assumes *mono f*

assumes $k < \text{size } M$

shows $f (\text{nth-mset } k M) = \text{nth-mset } k (\text{image-mset } f M)$

<proof>

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** $\text{card } (\text{least } (k+1) (\text{set-mset } A)) = k+1$

<proof>

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 $(\lambda x. \text{real } (\text{nat } (f x))) \in O[F](g)$
<proof>

lemma *landau-ceil*:

assumes $(\lambda \cdot. 1) \in O[F^\uparrow](g)$
assumes $f \in O[F^\uparrow](g)$
shows $(\lambda x. \text{real-of-int } [f x]) \in O[F^\uparrow](g)$
<proof>

lemma *landau-rat-ceil*:

assumes $(\lambda \cdot. 1) \in O[F^\uparrow](g)$
assumes $(\lambda x. \text{real-of-rat } (f x)) \in O[F^\uparrow](g)$
shows $(\lambda x. \text{real-of-int } [f x]) \in O[F^\uparrow](g)$
<proof>

lemma *landau-nat-ceil*:

assumes $(\lambda \cdot. 1) \in O[F^\uparrow](g)$
assumes $f \in O[F^\uparrow](g)$
shows $(\lambda x. \text{real } (\text{nat } [f x])) \in O[F^\uparrow](g)$
<proof>

lemma *eventually-prod1'*:

assumes $B \neq \text{bot}$
assumes $(\forall_F x \text{ in } A. P x)$
shows $(\forall_F x \text{ in } A \times_F B. P (\text{fst } x))$
<proof>

lemma *eventually-prod2'*:

assumes $A \neq \text{bot}$
assumes $(\forall_F x \text{ in } B. P x)$
shows $(\forall_F x \text{ in } A \times_F B. P (\text{snd } x))$
<proof>

lemma *sequentially-inf*: $\forall_F x \text{ in sequentially. } n \leq \text{real } x$
<proof>

instantiation *rat* :: *linorder-topology*
begin

definition *open-rat* :: *rat set* \Rightarrow *bool*

where *open-rat* = *generate-topology* ($\text{range } (\lambda a. \{.. < a\}) \cup \text{range } (\lambda a. \{a <..\})$)

instance

<proof>

end

lemma *inv-at-right-0-inf*:

$\forall_F x \text{ in at-right } 0. c \leq 1 / \text{real-of-rat } x$
<proof>

end

5 Probability Spaces

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

theory *Probability-Ext*

imports

HOL-Probability.Stream-Space

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

Frequency-Moments-Preliminary-Results

begin

Random variables that depend on disjoint sets of the components of a product space are independent.

lemma *make-ext*:

assumes $\bigwedge x. P\ x = P\ (\text{restrict } x\ I)$

shows $(\forall x \in \text{Pi } I\ A. P\ x) = (\forall x \in \text{PiE } I\ A. P\ x)$

<proof>

lemma *PiE-reindex*:

assumes *inj-on* $f\ I$

shows $\text{PiE } I\ (A \circ f) = (\lambda a. \text{restrict } (a \circ f)\ I)\ \text{PiE } (f\ \text{' } I)\ A$ (**is** *?lhs* = *?g* *'* *?rhs*)

<proof>

context *prob-space*

begin

lemma *indep-sets-reindex*:

assumes *inj-on* $f\ I$

shows $\text{indep-sets } A\ (f\ \text{' } I) = \text{indep-sets } (\lambda i. A\ (f\ i))\ I$

<proof>

lemma *indep-vars-cong-AE*:

assumes *AE* $x\ \text{in } M. (\forall i \in I. X'\ i\ x = Y'\ i\ x)$

assumes *indep-vars* $M'\ X'\ I$

assumes $\bigwedge i. i \in I \implies \text{random-variable } (M'\ i)\ (Y'\ i)$

shows *indep-vars* $M'\ Y'\ I$

<proof>

lemma *indep-vars-reindex*:

assumes *inj-on* $f\ I$

assumes *indep-vars* $M'\ X'\ (f\ \text{' } I)$

shows *indep-vars* $(M' \circ f)\ (\lambda k\ \omega. X'\ (f\ k)\ \omega)\ I$

<proof>

lemma *variance-divide*:

fixes $f :: 'a \Rightarrow \text{real}$
assumes $\text{integrable } M f$
shows $\text{variance } (\lambda\omega. f \ \omega / r) = \text{variance } f / r^2$
 $\langle \text{proof} \rangle$

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 $\text{prob } P \leq \text{prob } Q$
 $\langle \text{proof} \rangle$

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 $\text{prob } P \leq \text{prob } Q + \text{prob } R$
 $\langle \text{proof} \rangle$

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** $?lhs \leq ?rhs$)
 $\langle \text{proof} \rangle$

definition covariance where
 $\text{covariance } f \ g = \text{expectation } (\lambda\omega. (f \ \omega - \text{expectation } f) * (g \ \omega - \text{expectation } g))$

lemma real-prod-integrable:
fixes $f \ g :: 'a \Rightarrow \text{real}$
assumes $[\text{measurable}]: f \in \text{borel-measurable } M \ g \in \text{borel-measurable } M$
assumes $\text{sq-int}: \text{integrable } M (\lambda\omega. f \ \omega^2) \ \text{integrable } M (\lambda\omega. g \ \omega^2)$
shows $\text{integrable } M (\lambda\omega. f \ \omega * g \ \omega)$
 $\langle \text{proof} \rangle$

lemma covariance-eq:
fixes $f \ g :: 'a \Rightarrow \text{real}$
assumes $f \in \text{borel-measurable } M \ g \in \text{borel-measurable } M$
assumes $\text{integrable } M (\lambda\omega. f \ \omega^2) \ \text{integrable } M (\lambda\omega. g \ \omega^2)$
shows $\text{covariance } f \ g = \text{expectation } (\lambda\omega. f \ \omega * g \ \omega) - \text{expectation } f * \text{expectation } g$
 $\langle \text{proof} \rangle$

lemma covar-integrable:
fixes $f \ g :: 'a \Rightarrow \text{real}$
assumes $f \in \text{borel-measurable } M \ g \in \text{borel-measurable } M$
assumes $\text{integrable } M (\lambda\omega. f \ \omega^2) \ \text{integrable } M (\lambda\omega. g \ \omega^2)$
shows $\text{integrable } M (\lambda\omega. (f \ \omega - \text{expectation } f) * (g \ \omega - \text{expectation } g))$
 $\langle \text{proof} \rangle$

lemma *sum-square-int*:

fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$

assumes *finite I*

assumes $\bigwedge i. i \in I \implies f\ i \in \text{borel-measurable } M$

assumes $\bigwedge i. i \in I \implies \text{integrable } M (\lambda\omega. f\ i\ \omega^{\wedge}2)$

shows *integrable* $M (\lambda\omega. (\sum i \in I. f\ i\ \omega)^2)$

<proof>

lemma *var-sum-1*:

fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$

assumes *finite I*

assumes $\bigwedge i. i \in I \implies f\ i \in \text{borel-measurable } M$

assumes $\bigwedge i. i \in I \implies \text{integrable } M (\lambda\omega. f\ i\ \omega^{\wedge}2)$

shows

variance $(\lambda\omega. (\sum i \in I. f\ i\ \omega)) = (\sum i \in I. (\sum j \in I. \text{covariance } (f\ i) (f\ j)))$

<proof>

lemma *covar-self-eq*:

fixes $f :: 'a \Rightarrow \text{real}$

shows *covariance* $f\ f = \text{variance } f$

<proof>

lemma *covar-indep-eq-zero*:

fixes $f\ g :: 'a \Rightarrow \text{real}$

assumes *integrable* $M\ f$

assumes *integrable* $M\ g$

assumes *indep-var borel* $f\ \text{borel } g$

shows *covariance* $f\ g = 0$

<proof>

lemma *var-sum-2*:

fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$

assumes *finite I*

assumes $\bigwedge i. i \in I \implies f\ i \in \text{borel-measurable } M$

assumes $\bigwedge i. i \in I \implies \text{integrable } M (\lambda\omega. f\ i\ \omega^{\wedge}2)$

shows *variance* $(\lambda\omega. (\sum i \in I. f\ i\ \omega)) =$

$(\sum i \in I. \text{variance } (f\ i)) + (\sum i \in I. \sum j \in I - \{i\}. \text{covariance } (f\ i) (f\ j))$

<proof>

lemma *var-sum-pairwise-indep*:

fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$

assumes *finite I*

assumes $\bigwedge i. i \in I \implies f\ i \in \text{borel-measurable } M$

assumes $\bigwedge i. i \in I \implies \text{integrable } M (\lambda\omega. f\ i\ \omega^{\wedge}2)$

assumes $\bigwedge i\ j. i \in I \implies j \in I \implies i \neq j \implies \text{indep-var borel } (f\ i)\ \text{borel } (f\ j)$

shows *variance* $(\lambda\omega. (\sum i \in I. f\ i\ \omega)) = (\sum i \in I. \text{variance } (f\ i))$

<proof>

lemma *indep-var-from-indep-vars*:

```

assumes  $i \neq j$ 
assumes indep-vars ( $\lambda \cdot. M'$ )  $f \{i, j\}$ 
shows indep-var  $M' (f i) M' (f j)$ 
<proof>

```

lemma *var-sum-pairwise-indep-2*:

```

fixes  $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$ 
assumes finite  $I$ 
assumes  $\bigwedge i. i \in I \implies f i \in \text{borel-measurable } M$ 
assumes  $\bigwedge i. i \in I \implies \text{integrable } M (\lambda \omega. f i \omega^{\mathcal{Q}})$ 
assumes  $\bigwedge J. J \subseteq I \implies \text{card } J = 2 \implies \text{indep-vars } (\lambda \cdot. \text{borel}) f J$ 
shows variance ( $\lambda \omega. (\sum i \in I. f i \omega)$ ) = ( $\sum i \in I. \text{variance } (f i)$ )
<proof>

```

lemma *var-sum-all-indep*:

```

fixes  $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$ 
assumes finite  $I$ 
assumes  $\bigwedge i. i \in I \implies f i \in \text{borel-measurable } M$ 
assumes  $\bigwedge i. i \in I \implies \text{integrable } M (\lambda \omega. f i \omega^{\mathcal{Q}})$ 
assumes indep-vars ( $\lambda \cdot. \text{borel}$ )  $f I$ 
shows variance ( $\lambda \omega. (\sum i \in I. f i \omega)$ ) = ( $\sum i \in I. \text{variance } (f i)$ )
<proof>

```

end

end

6 Indexed Products of Probability Mass Functions

theory *Product-PMF-Ext*

imports *Main Probability-Ext HOL-Probability.Product-PMF Universal-Hash-Families.Preliminary-Results*
begin

hide-const *Isolated.discrete*

This section introduces a restricted version of *Pi-pmf* where the default value is *undefined* and contains some additional results about that case in addition to *HOL-Probability.Product-PMF*

abbreviation *prod-pmf* **where** *prod-pmf* $I M \equiv \text{Pi-pmf } I \text{ undefined } M$

lemma *pmf-prod-pmf*:

```

assumes finite  $I$ 
shows pmf (prod-pmf  $I M$ )  $x = (\text{if } x \in \text{extensional } I \text{ then } \prod i \in I. (\text{pmf } (M i)) (x i) \text{ else } 0)$ 
<proof>

```

lemma *PiE-default-undefined-eq*: *PiE-dflt* $I \text{ undefined } M = \text{PiE } I M$

<proof>

lemma *set-prod-pmf*:

assumes *finite I*

shows $set\text{-}pmf (prod\text{-}pmf I M) = PiE I (set\text{-}pmf \circ M)$

<proof>

A more general version of *measure-Pi-pmf-Pi*.

lemma *prob-prod-pmf'*:

assumes *finite I*

assumes $J \subseteq I$

shows $measure (measure\text{-}pmf (Pi\text{-}pmf I d M)) (Pi J A) = (\prod_{i \in J} measure (M i) (A i))$

<proof>

lemma *prob-prod-pmf-slice*:

assumes *finite I*

assumes $i \in I$

shows $measure (measure\text{-}pmf (prod\text{-}pmf I M)) \{\omega. P (\omega i)\} = measure (M i) \{\omega. P \omega\}$

<proof>

definition *restrict-dfl* **where** $restrict\text{-}dfl f A d = (\lambda x. if x \in A then f x else d)$

lemma *pi-pmf-decompose*:

assumes *finite I*

shows $Pi\text{-}pmf I d M = map\text{-}pmf (\lambda \omega. restrict\text{-}dfl (\lambda i. \omega (f i) i) I d) (Pi\text{-}pmf (f ' I) (\lambda -. d) (\lambda j. Pi\text{-}pmf (f -' \{j\} \cap I) d M))$

<proof>

lemma *restrict-dfl-iter*: $restrict\text{-}dfl (restrict\text{-}dfl f I d) J d = restrict\text{-}dfl f (I \cap J) d$

<proof>

lemma *indep-vars-restrict'*:

assumes *finite I*

shows $prob\text{-}space.indep\text{-}vars (Pi\text{-}pmf I d M) (\lambda -. discrete) (\lambda i \omega. restrict\text{-}dfl \omega (f -' \{i\} \cap I) d) (f ' I)$

<proof>

lemma *indep-vars-restrict-intro'*:

assumes *finite I*

assumes $\bigwedge i \omega. i \in J \implies X' i \omega = X' i (restrict\text{-}dfl \omega (f -' \{i\} \cap I) d)$

assumes $J = f ' I$

assumes $\bigwedge \omega i. i \in J \implies X' i \omega \in space (M' i)$

shows $prob\text{-}space.indep\text{-}vars (measure\text{-}pmf (Pi\text{-}pmf I d p)) M' (\lambda i \omega. X' i \omega) J$

<proof>

lemma

```

fixes  $f :: 'b \Rightarrow ('c :: \{second-countable-topology,banach,real-normed-field\})$ 
assumes  $finite\ I$ 
assumes  $i \in I$ 
assumes  $integrable\ (measure-pmf\ (M\ i))\ f$ 
shows  $integrable-Pi-pmf-slice: integrable\ (Pi-pmf\ I\ d\ M)\ (\lambda x. f\ (x\ i))$ 
and  $expectation-Pi-pmf-slice: integral^L\ (Pi-pmf\ I\ d\ M)\ (\lambda x. f\ (x\ i)) = integral^L$ 
 $(M\ i)\ f$ 
 $\langle proof \rangle$ 

```

This is an improved version of *expectation-prod-Pi-pmf*. It works for general normed fields instead of non-negative real functions .

```

lemma  $expectation-prod-Pi-pmf:$ 
fixes  $f :: 'a \Rightarrow 'b \Rightarrow ('c :: \{second-countable-topology,banach,real-normed-field\})$ 
assumes  $finite\ I$ 
assumes  $\bigwedge i. i \in I \implies integrable\ (measure-pmf\ (M\ i))\ (f\ i)$ 
shows  $integral^L\ (Pi-pmf\ I\ d\ M)\ (\lambda x. (\prod i \in I. f\ i\ (x\ i))) = (\prod i \in I. integral^L$ 
 $(M\ i)\ (f\ i))$ 
 $\langle proof \rangle$ 

```

```

lemma  $variance-prod-pmf-slice:$ 
fixes  $f :: 'a \Rightarrow real$ 
assumes  $i \in I\ finite\ I$ 
assumes  $integrable\ (measure-pmf\ (M\ i))\ (\lambda \omega. f\ \omega^{\wedge} 2)$ 
shows  $prob-space.variance\ (Pi-pmf\ I\ d\ M)\ (\lambda \omega. f\ (\omega\ i)) = prob-space.variance$ 
 $(M\ i)\ f$ 
 $\langle proof \rangle$ 

```

```

lemma  $Pi-pmf-bind-return:$ 
assumes  $finite\ I$ 
shows  $Pi-pmf\ I\ d\ (\lambda i. M\ i \gg (\lambda x. return-pmf\ (f\ i\ x))) = Pi-pmf\ I\ d'\ M \gg$ 
 $(\lambda x. return-pmf\ (\lambda i. if\ i \in I\ then\ f\ i\ (x\ i)\ else\ d))$ 
 $\langle proof \rangle$ 

```

end

7 Frequency Moment 0

```

theory  $Frequency-Moment-0$ 
imports
   $Frequency-Moments-Preliminary-Results$ 
   $Median-Method.Median$ 
   $K-Smallest$ 
   $Universal-Hash-Families.Carter-Wegman-Hash-Family$ 
   $Frequency-Moments$ 
   $Landau-Ext$ 
   $Product-PMF-Ext$ 
   $Universal-Hash-Families.Field$ 
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 $r = \text{nat } (4 * \lceil \log 2 (1 / \text{real-of-rat } \delta) \rceil + 23)$;
 $h \leftarrow \text{prod-pmf } \{..\langle s \rangle\} (\lambda\cdot. \text{pmf-of-set } (\text{bounded-degree-polynomials } (\text{mod-ring } p) \ 2))$;
return-pmf $(s, t, p, r, h, (\lambda\cdot \in \{0..\langle s \rangle\}. \{\}))$
}

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 $(s, t, p, r, h, \lambda i \in \{..\langle s \rangle\}.$
least t $(\text{insert } (\text{float-of } (\text{truncate-down } r (\text{hash } p \ x \ (h \ i)))) (\text{sketch } i)))$

fun $f0\text{-result} :: f0\text{-state} \Rightarrow \text{rat pmf where}$
 $f0\text{-result } (s, t, p, r, h, \text{sketch}) = \text{return-pmf } (\text{median } s (\lambda i \in \{..\langle s \rangle\}.$
(if $\text{card } (\text{sketch } i) < t$ then $\text{of-nat } (\text{card } (\text{sketch } i))$ 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) = ($$
let $s = \text{nat } \lceil -18 * \ln (\text{real-of-rat } \varepsilon) \rceil$ in
let $r = \text{nat } (4 * \lceil \log 2 (1 / \text{real-of-rat } \delta) \rceil + 23)$ in
let $t = \text{nat } \lceil 80 / (\text{real-of-rat } \delta)^2 \rceil$ in
 $6 +$
 $2 * \log 2 (\text{real } s + 1) +$
 $2 * \log 2 (\text{real } t + 1) +$
 $2 * \log 2 (\text{real } n + 21) +$
 $2 * \log 2 (\text{real } r + 1) +$
 $\text{real } s * (5 + 2 * \log 2 (21 + \text{real } n)) +$
 $\text{real } t * (13 + 4 * r + 2 * \log 2 (\log 2 (\text{real } n + 13))))$$

definition $\text{encode-f0-state} :: f0\text{-state} \Rightarrow \text{bool list option where}$

```

encode-f0-state =
  N_e ⋈_e (λs.
    N_e ×_e (
      N_e ⋈_e (λp.
        N_e ×_e (
          ([0..<s] →_e (P_e p 2)) ×_e
          ([0..<s] →_e (S_e F_e))))))

```

lemma *inj-on encode-f0-state* (dom encode-f0-state)
 ⟨proof⟩

context

```

fixes ε δ :: rat
fixes n :: nat
fixes as :: nat list
fixes result
assumes ε-range: ε ∈ {0<..assumes δ-range: δ ∈ {0<..assumes as-range: set as ⊆ {..defines result ≡ fold (λa state. state ≫= f0-update a) as (f0-init δ ε n) ≫=
f0-result
begin

```

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

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

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

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

```

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

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

```

  ⟨proof⟩ lemma r-ge-23:  $r \geq 23$ 
  ⟨proof⟩ lemma two-pow-r-le-1:  $0 < 1 - 2^{\text{powr } - \text{real } r}$ 
  ⟨proof⟩

```

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

```

rewrites ring.hash (mod-ring  $p$ ) = Frequency-Moment-0.hash  $p$ 
  ⟨proof⟩ definition tr-hash where  $\text{tr-hash } x \ \omega = \text{truncate-down } r (\text{hash } x \ \omega)$ 

```

private definition *sketch-rv where*

```

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

```

private definition *estimate*

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

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

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

private definition Ω_0 **where** $\Omega_0 = \text{prod-pmf } \{..<s\} (\lambda-. \text{pmf-of-set space})$

private lemma *f0-alg-sketch*:

defines *sketch* $\equiv \text{fold } (\lambda a \text{ state. state } \gg= \text{f0-update } a) \text{ as (f0-init } \delta \ \varepsilon \ n)$

shows *sketch* $= \text{map-pmf } (\lambda x. (s, t, p, r, x, \lambda i \in \{..<s\}. \text{sketch-rv } (x \ i))) \ \Omega_0$

<proof> **lemma** *card-nat-in-ball*:

fixes $x :: \text{nat}$

fixes $q :: \text{real}$

assumes $q \geq 0$

defines $A \equiv \{k. \text{abs (real } x - \text{real } k) \leq q \wedge k \neq x\}$

shows $\text{real (card } A) \leq 2 * q \text{ and finite } A$

<proof> **lemma** *prob-degree-lt-1*:

$\text{prob } \{\omega. \text{degree } \omega < 1\} \leq 1 / \text{real } p$

<proof> **lemma** *collision-prob*:

assumes $c \geq 1$

shows $\text{prob } \{\omega. \exists x \in \text{set as. } \exists y \in \text{set as. } x \neq y \wedge \text{tr-hash } x \ \omega \leq c \wedge \text{tr-hash } x \ \omega = \text{tr-hash } y \ \omega\} \leq$

$(5/2) * (\text{real (card (set as))})^2 * c^2 * 2 \text{ powr } -(\text{real } r) / (\text{real } p)^2 + 1 / \text{real } p$

(is prob $\{\omega. ?l \ \omega\} \leq ?r1 + ?r2)$

<proof> **lemma** *of-bool-square*: $(\text{of-bool } x)^2 = ((\text{of-bool } x)::\text{real})$

<proof> **definition** Q **where** $Q \ y \ \omega = \text{card } \{x \in \text{set as. int (hash } x \ \omega) < y\}$

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

private lemma

assumes $a \geq 0$

assumes $a \leq \text{int } p$

shows *exp-Q*: $\text{expectation } (\lambda \omega. \text{real } (Q \ a \ \omega)) = \text{real } m * (\text{of-int } a) / p$

and *var-Q*: $\text{variance } (\lambda \omega. \text{real } (Q \ a \ \omega)) \leq \text{real } m * (\text{of-int } a) / p$

<proof> **lemma** *t-bound*: $t \leq 81 / (\text{real-of-rat } \delta)^2$

<proof> **lemma** *t-r-bound*:

$18 * 40 * (\text{real } t)^2 * 2 \text{ powr } (-\text{real } r) \leq 1$

<proof> **lemma** *m-eq-F-0*: $\text{real } m = \text{of-rat } (F \ 0 \ \text{as})$

<proof> **lemma** *estimate'-bounds*:

$\text{prob } \{\omega. \text{of-rat } \delta * \text{real-of-rat } (F \ 0 \ \text{as}) < |\text{estimate}' (\text{sketch-rv}' \ \omega) - \text{of-rat } (F \ 0 \ \text{as})|\} \leq 1/3$

<proof> **lemma** *median-bounds*:

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

<proof>

lemma *f0-alg-correct'*:

$\mathcal{P}(\omega \text{ in measure-pmf result. } |\omega - F \ 0 \ as| \leq \delta * F \ 0 \ as) \geq 1 - \text{of-rat } \varepsilon$

<proof> **lemma** *f-subset*:

assumes $g \ ' \ A \subseteq h \ ' \ B$

shows $(\lambda x. f \ (g \ x)) \ ' \ A \subseteq (\lambda x. f \ (h \ x)) \ ' \ B$

<proof>

lemma *f0-exact-space-usage'*:

defines $\Omega \equiv \text{fold } (\lambda a \ \text{state. state} \gg= \text{f0-update } a) \ \text{as } (\text{f0-init } \delta \ \varepsilon \ n)$

shows $AE \ \omega \ \text{in } \Omega. \text{bit-count } (\text{encode-f0-state } \omega) \leq \text{f0-space-usage } (n, \varepsilon, \delta)$

<proof>

end

Main results of this section:

theorem *f0-alg-correct*:

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

assumes $\delta \in \{0 < .. < 1\}$

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

defines $\Omega \equiv \text{fold } (\lambda a \ \text{state. state} \gg= \text{f0-update } a) \ \text{as } (\text{f0-init } \delta \ \varepsilon \ n) \gg= \text{f0-result}$

shows $\mathcal{P}(\omega \text{ in measure-pmf } \Omega. |\omega - F \ 0 \ as| \leq \delta * F \ 0 \ as) \geq 1 - \text{of-rat } \varepsilon$

<proof>

theorem *f0-exact-space-usage*:

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

assumes $\delta \in \{0 < .. < 1\}$

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

defines $\Omega \equiv \text{fold } (\lambda a \ \text{state. state} \gg= \text{f0-update } a) \ \text{as } (\text{f0-init } \delta \ \varepsilon \ n)$

shows $AE \ \omega \ \text{in } \Omega. \text{bit-count } (\text{encode-f0-state } \omega) \leq \text{f0-space-usage } (n, \varepsilon, \delta)$

<proof>

theorem *f0-asymptotic-space-complexity*:

$\text{f0-space-usage} \in O[\text{at-top } \times_F \ \text{at-right } 0 \ \times_F \ \text{at-right } 0](\lambda(n, \varepsilon, \delta). \ln(1 / \text{of-rat } \varepsilon) *$

$(\ln(\text{real } n) + 1 / (\text{of-rat } \delta)^2 * (\ln(\ln(\text{real } n)) + \ln(1 / \text{of-rat } \delta))))$

$(\text{is } - \in O[?F](?rhs))$

<proof>

end

8 Frequency Moment 2

theory *Frequency-Moment-2*

imports

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

Equivalence-Relation-Enumeration.Equivalence-Relation-Enumeration

Landau-Ext

Median-Method.Median
Product-PMF-Ext
Universal-Hash-Families.Field
Frequency-Moments

begin

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 $p\ h\ k = (\text{if even } (\text{ring.hash } (\text{mod-ring } p)\ k\ h) \text{ then int } p - 1 \text{ else } - \text{int } p - 1)$

type-synonym *f2-state* = $\text{nat} \times \text{nat} \times \text{nat} \times (\text{nat} \times \text{nat} \Rightarrow \text{nat list}) \times (\text{nat} \times \text{nat} \Rightarrow \text{int})$

fun *f2-init* :: $\text{rat} \Rightarrow \text{rat} \Rightarrow \text{nat} \Rightarrow \text{f2-state pmf}$ **where**

f2-init $\delta\ \varepsilon\ n =$
do {
let $s_1 = \text{nat } \lceil 6 / \delta^2 \rceil$;
let $s_2 = \text{nat } \lceil -(18 * \ln (\text{real-of-rat } \varepsilon)) \rceil$;
let $p = \text{prime-above } (\text{max } n\ 3)$;
 $h \leftarrow \text{prod-pmf } (\{\dots < s_1\} \times \{\dots < s_2\}) (\lambda\cdot. \text{pmf-of-set } (\text{bounded-degree-polynomials } (\text{mod-ring } p)\ 4))$;
return-pmf $(s_1, s_2, p, h, (\lambda\ \in \{\dots < s_1\} \times \{\dots < s_2\}. (0 :: \text{int})))$
}

fun *f2-update* :: $\text{nat} \Rightarrow \text{f2-state} \Rightarrow \text{f2-state pmf}$ **where**

f2-update $x\ (s_1, s_2, p, h, \text{sketch}) =$
return-pmf $(s_1, s_2, p, h, \lambda i \in \{\dots < s_1\} \times \{\dots < s_2\}. \text{f2-hash } p\ (h\ i)\ x + \text{sketch } i)$

fun *f2-result* :: $\text{f2-state} \Rightarrow \text{rat pmf}$ **where**

f2-result $(s_1, s_2, p, h, \text{sketch}) =$
return-pmf $(\text{median } s_2\ (\lambda i_2 \in \{\dots < s_2\}. (\sum_{i_1 \in \{\dots < s_1\}} . (\text{rat-of-int } (\text{sketch } (i_1, i_2))))^2) / (((\text{rat-of-nat } p)^2 - 1) * \text{rat-of-nat } s_1))$

fun *f2-space-usage* :: $(\text{nat} \times \text{nat} \times \text{rat} \times \text{rat}) \Rightarrow \text{real}$ **where**

f2-space-usage $(n, m, \varepsilon, \delta) = ($
let $s_1 = \text{nat } \lceil 6 / \delta^2 \rceil$ in
let $s_2 = \text{nat } \lceil -(18 * \ln (\text{real-of-rat } \varepsilon)) \rceil$ in
 $3 +$
 $2 * \log 2\ (s_1 + 1) +$
 $2 * \log 2\ (s_2 + 1) +$
 $2 * \log 2\ (9 + 2 * \text{real } n) +$
 $s_1 * s_2 * (5 + 4 * \log 2\ (8 + 2 * \text{real } n) + 2 * \log 2\ (\text{real } m * (18 + 4 * \text{real } n) + 1))$
 $)$

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

```

encode-f2-state =
   $N_e \times_e (\lambda s_1.$ 
     $N_e \times_e (\lambda s_2.$ 
       $N_e \times_e (\lambda p.$ 
        (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*)

<proof>

context

fixes $\varepsilon \delta :: \text{rat}$

fixes $n :: \text{nat}$

fixes *as* :: *nat list*

fixes *result*

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

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

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

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

begin

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

lemma *s1-gt-0*: $s_1 > 0$

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

lemma *s2-gt-0*: $s_2 > 0$

<proof> **definition** p **where** $p = \text{prime-above } (\text{max } n \ 3)$

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

<proof>

lemma *p-ge-3*: $p \geq 3$

<proof>

lemma *p-gt-0*: $p > 0$ *<proof>*

lemma *p-gt-1*: $p > 1$ *<proof>*

lemma *p-ge-n*: $p \geq n$ *<proof>*

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

<proof>

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

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

space)

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

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

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

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

lemma *mean-rv-alg-sketch*:

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

lemma *distr*: *result* = *map-pmf result-rv* Ω

<proof> **lemma** *f2-hash-pow-exp*:

assumes $k < p$

shows

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

lemma

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

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

<proof>

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

<proof>

interpretation Ω : *prob-space* Ω_p

<proof>

lemma *integrable- Ω* :

fixes $f :: ((\text{nat} \times \text{nat}) \Rightarrow (\text{nat list})) \Rightarrow \text{real}$

shows *integrable* $\Omega_p \ f$

<proof>

lemma *sketch-rv-exp*:

assumes $i_2 < s_2$

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

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

<proof>

lemma *sketch-rv-var*:

assumes $i_2 < s_2$
assumes $i_1 \in \{0..<s_1\}$
shows $\Omega.\text{variance } (\lambda\omega. \text{sketch-rv } (\omega (i_1, i_2))) \leq 2 * (\text{real-of-rat } (F \ 2 \ as))^2 * ((\text{real } p)^2 - 1)^2$
 <proof>

lemma *mean-rv-exp*:

assumes $i < s_2$
shows $\Omega.\text{expectation } (\lambda\omega. \text{mean-rv } \omega \ i) = \text{real-of-rat } (F \ 2 \ as)$
 <proof>

lemma *mean-rv-var*:

assumes $i < s_2$
shows $\Omega.\text{variance } (\lambda\omega. \text{mean-rv } \omega \ i) \leq (\text{real-of-rat } (\delta * F \ 2 \ as))^2 / 3$
 <proof>

lemma *mean-rv-bounds*:

assumes $i < s_2$
shows $\Omega.\text{prob } \{\omega. \text{real-of-rat } \delta * \text{real-of-rat } (F \ 2 \ as) < |\text{mean-rv } \omega \ i - \text{real-of-rat } (F \ 2 \ as)|\} \leq 1/3$
 <proof>

lemma *f2-alg-correct'*:

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

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

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

end

Main results of this section:

theorem *f2-alg-correct*:

assumes $\varepsilon \in \{0<..
assumes $\delta > 0$
assumes $\text{set } as \subseteq \{..
defines $\Omega \equiv \text{fold } (\lambda a \text{ state. state } \ggg \text{f2-update } a) \ as \ (\text{f2-init } \delta \ \varepsilon \ n) \ \ggg \ \text{f2-result}$
shows $\mathcal{P}(\omega \text{ in measure-pmf } \Omega. |\omega - F \ 2 \ as| \leq \delta * F \ 2 \ as) \geq 1 - \text{of-rat } \varepsilon$
 <proof>$$

theorem *f2-exact-space-usage*:

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

<proof>

theorem *f2-asymptotic-space-complexity*:

f2-space-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

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-state* = $nat \times nat \times nat \times nat \times (nat \times nat \Rightarrow (nat \times nat))$

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

fk-init $k\ \delta\ \varepsilon\ n =$

do {

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

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

return-*pmf* ($s_1, s_2, k, 0, (\lambda - \in \{0..<s_1\} \times \{0..<s_2\}. (0,0))$)

}

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

fk-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-*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 $(x, l + of\text{-}bool (x = a))$

)

)

}

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

fk-result (s_1, s_2, k, m, r) =

```

return-pmf (median s2 (λi2 ∈ {0..<s2}.
  (∑ i1 ∈ {0..<s1}. rat-of-nat (let t = snd (r (i1, i2)) + 1 in m * (t^k - (t -
1) ^k))) / (rat-of-nat s1))
)

```

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

```

fun fk-space-usage :: (nat × nat × nat × rat × rat) ⇒ real where
fk-space-usage (k, n, m, ε, δ) = (
  let s1 = nat ⌈3*real k* (real n) powr (1-1/ real k) / (real-of-rat δ)2 ⌋ in
  let s2 = nat ⌈-(18 * ln (real-of-rat ε))⌋ in
  4 +
  2 * log 2 (s1 + 1) +
  2 * log 2 (s2 + 1) +
  2 * log 2 (real k + 1) +
  2 * log 2 (real m + 1) +
  s1 * s2 * (2 + 2 * log 2 (real n+1) + 2 * log 2 (real m+1)))

```

definition *encode-fk-state* :: *fk-state* ⇒ *bool list option* **where**
encode-fk-state =
 $N_e \times_e (\lambda s_1.$
 $N_e \times_e (\lambda s_2.$
 $N_e \times_e$
 $N_e \times_e$
 $(List.product [0..<s_1] [0..<s_2] \rightarrow_e (N_e \times_e N_e)))$)

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 ⇒ (nat × 'a × nat) ⇒ (nat × 'a × nat) pmf where
fk-update-2 a (m,x,l) =
  do {
    coin ← 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 $fold (\lambda x s. s \gg= fk-update-2 x) as (return-pmf (0,0,0)) =$
 $pmf-of-set \{..<length as\} \gg= (\lambda k. return-pmf (length as, sketch as k))$
<proof>

context
fixes $\varepsilon \delta :: rat$

fixes $n\ k :: \text{nat}$
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 $as\text{-range}: \text{set } as \subseteq \{.. < n\}$
begin

definition s_1 **where** $s_1 = \text{nat } \lceil \beta * \text{real } k * (\text{real } n) \text{ powr } (1 - 1 / \text{real } k) / (\text{real-of-rat } \delta)^2 \rceil$
definition s_2 **where** $s_2 = \text{nat } \lceil -(18 * \ln (\text{real-of-rat } \varepsilon)) \rceil$

definition $M_1 = \{(u, v). v < \text{count-list } as\ u\}$
definition $\Omega_1 = \text{measure-pmf } (\text{pmf-of-set } M_1)$

definition $M_2 = \text{prod-pmf } (\{0.. < s_1\} \times \{0.. < s_2\}) (\lambda-. \text{pmf-of-set } M_1)$
definition $\Omega_2 = \text{measure-pmf } M_2$

interpretation $\text{prob-space } \Omega_1$
 $\langle \text{proof} \rangle$

interpretation $\Omega_2: \text{prob-space } \Omega_2$
 $\langle \text{proof} \rangle$

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

lemma
assumes $as \neq []$
shows $\text{fin-space}: \text{finite } M_1$
and $\text{non-empty-space}: M_1 \neq \{\}$
and $\text{card-space}: \text{card } M_1 = \text{length } as$
 $\langle \text{proof} \rangle$

lemma
assumes $as \neq []$
shows $\text{integrable-1}: \text{integrable } \Omega_1 (f :: - \Rightarrow \text{real})$ **and**
 $\text{integrable-2}: \text{integrable } \Omega_2 (g :: - \Rightarrow \text{real})$
 $\langle \text{proof} \rangle$

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

lemma $\text{fk-update-distr}: \text{fold } (\lambda x\ s. s \ggg \text{fk-update } x) as (\text{fk-init } k\ \delta\ \varepsilon\ n) =$

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

lemma *power-diff-sum*:

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

lemma *power-diff-est*:

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

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 $|\sum x \in A. f \ x * g \ x| \leq (\sum x \in A. |f \ x| \text{ powr } p) \text{ powr } (1/p) * (\sum x \in A. |g \ x| \text{ powr } q) \text{ powr } (1/q)$
 <proof>

lemma *real-count-list-pos*:

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

lemma *fk-estimate*:

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

definition *result*

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

lemma *result-exp-1*:

assumes $\text{as} \neq []$
shows $\text{expectation result} = \text{real-of-rat } (F \ k \ \text{as})$
 <proof>

lemma *result-var-1*:

assumes $\text{as} \neq []$

shows $\text{variance result} \leq (\text{of-rat } (F k as))^2 * k * n \text{ powr } (1 - 1 / \text{real } k)$
 ⟨proof⟩

theorem *fk-alg-sketch*:

assumes $as \neq []$

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

⟨proof⟩

definition *mean-rv*

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

definition *median-rv*

where $\text{median-rv } \omega = \text{median } s_2 (\lambda i_2. \text{mean-rv } \omega \ i_2)$

lemma *fk-alg-correct'*:

defines $M \equiv \text{fold } (\lambda a \text{ state. state} \ggg \text{fk-update } a) \text{ as } (\text{fk-init } k \delta \varepsilon n) \ggg \text{fk-result}$

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

⟨proof⟩

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

defines $M \equiv \text{fold } (\lambda a \text{ state. state} \ggg \text{fk-update } a) \text{ as } (\text{fk-init } k \delta \varepsilon n)$

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

(is $AE \ \omega \text{ in } M. (- \leq \text{?rhs})$)

⟨proof⟩

end

Main results of this section:

theorem *fk-alg-correct*:

assumes $k \geq 1$

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

assumes $\delta > 0$

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

defines $M \equiv \text{fold } (\lambda a \text{ state. state} \ggg \text{fk-update } a) \text{ as } (\text{fk-init } k \delta \varepsilon n) \ggg \text{fk-result}$

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

⟨proof⟩

theorem *fk-exact-space-usage*:

assumes $k \geq 1$

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

assumes $\delta > 0$

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

defines $M \equiv \text{fold } (\lambda a \text{ state. state} \ggg \text{fk-update } a) \text{ as } (\text{fk-init } k \delta \varepsilon n)$

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

⟨proof⟩

theorem *fk-asymptotic-space-complexity:*

fk-space-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.¹ With these 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 $F_0 < t$.

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

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 $Q(u) = \sum_{a \in A} 1_{\{h(a) < u\}}$ ², i.e., Q is a sum of pairwise independent $\{0, 1\}$ -valued random variables, with expectation $\frac{u}{p}$ and variance $\frac{u}{p} - \frac{u^2}{p^2}$.

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

³ 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 $v = \lfloor \frac{tp}{(1-\delta')F_0} \rfloor$ it is possible to conclude:

$$t-1 \leq 4 \frac{t}{(1-\delta')} - 3\sqrt{\frac{t}{(1-\delta')}} - 1 \leq \frac{F_0v}{p} - 3\sqrt{\frac{F_0v}{p}} \leq \text{E}Q(v) - 3\sqrt{\text{Var}Q(v)}$$

and thus using Tchebyshev's inequality:

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

Similarly for $u = \lfloor \frac{tp}{(1+\delta')F_0} \rfloor$ 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:

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

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, 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 $1 - \frac{1}{p}$,

³A consequence of h being chosen uniformly from a 2-independent hash family.

⁴The verification of this inequality is a lengthy but straightforward calculation using the definition of δ' and t .

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:

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

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.