

Formalization of Randomized Approximation Algorithms for Frequency Moments

Emin Karayel

April 19, 2022

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	6
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 ∧ n M x . P M ⇒ ¬(x ∈ # M) ⇒ n > 0 ⇒ P (M + *replicate-mset* n x)
shows P M
 ⟨*proof*⟩

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

lemma *sum-collapse*:
fixes f :: 'a ⇒ 'b::{*comm-monoid-add*}
assumes *finite* A
assumes z ∈ A
assumes ∧ y . y ∈ A ⇒ y ≠ z ⇒ f y = 0
shows *sum* f A = f z
 ⟨*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 ⇒ 'b::{*ring,semiring-1*}
shows *sum-list* (*map* f xs) = (∑ x ∈ *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 $prod-list (map f xs) = (\prod x \in set xs. (f x) \wedge (count-list xs x))$
<proof>

lemma *sorted-sorted-list-of-multiset*: $sorted (sorted-list-of-multiset M)$

<proof>

lemma *count-mset*: $count (mset xs) a = count-list xs a$

<proof>

lemma *swap-filter-image*: $filter-mset g (image-mset f A) = image-mset f (filter-mset (g \circ f) A)$

<proof>

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*:

$x \leq round-down\ prec\ x + 2\ powr (-prec)$

<proof>

lemma *truncate-down-ge*:

$x \leq truncate-down\ prec\ x + abs\ x * 2\ powr (-prec)$

<proof>

lemma *truncate-down-pos*:

assumes $x \geq 0$

shows $x * (1 - 2\ powr (-prec)) \leq truncate-down\ prec\ x$

<proof>

lemma *truncate-down-eq*:

assumes $truncate-down\ r\ x = truncate-down\ r\ y$

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

<proof>

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

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

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

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

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

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*: $\text{prime } (\text{prime-above } n)$
 and *prime-above-range*: $\text{prime-above } n \in \{n..(2*n+2)\}$
 ⟨proof⟩

lemma *prime-above-min*: $\text{prime-above } n \geq 2$
 ⟨proof⟩

lemma *prime-above-lower-bound*: $\text{prime-above } n \geq n$
 ⟨proof⟩

lemma *prime-above-upper-bound*: $\text{prime-above } n \leq 2*n+2$
 ⟨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**

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

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

<proof>

lemma *F-gr-0*:

assumes $as \neq []$

shows $F\ k\ as > 0$

<proof>

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

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

lemma *poly-encoding*:

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

<proof>

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

assumes $p > 1$

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

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

<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* f T
assumes $x \in T$
shows $\text{rank-of } x \ S = \text{rank-of } (f \ x) \ (f \ ' \ S)$
 <proof>

definition *least* **where** *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* ($\lambda x. \text{rank-of } x \ S$) 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* f S

shows $f \cdot \text{least } k \ S = \text{least } k \ (f \cdot S)$
<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* $:: \text{nat} \implies ('a :: \text{linorder}) \text{multiset} \implies 'a$ **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 $nth\text{-mset } k A = Max (least (k+1) (set\text{-mset } A))$ **and** $card (least (k+1) (set\text{-mset } A)) = k+1$
 ⟨*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 $f :: 'a \Rightarrow \text{int}$
assumes $(\lambda x. \text{of-int } (f x)) \in O[F](g)$
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 $(\lambda x. \text{real-of-int } \lceil f x \rceil) \in O[F^\uparrow](g)$
<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 $(\lambda x. \text{real-of-int } \lceil f x \rceil) \in O[F^\uparrow](g)$
<proof>

lemma *landau-nat-ceil*:

assumes $(\lambda \cdot. 1) \in O[F^\uparrow](g)$
assumes $f \in O[F^\uparrow](g)$
shows $(\lambda x. \text{real } (\text{nat } \lceil f x \rceil)) \in O[F^\uparrow](g)$
<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*:
 $\forall_F x \text{ in } \text{at-right } 0. c \leq 1 / \text{real-of-rat } x$
<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 $\text{PiE } I (A \circ f) = (\lambda a. \text{restrict } (a \circ f) I) \circ \text{PiE } (f \circ I) A$ (*is ?lhs = ?g* *?*
rhs)

<proof>

context *prob-space*

begin

lemma *indep-sets-reindex*:

assumes *inj-on* $f I$

shows $\text{indep-sets } A (f \circ I) = \text{indep-sets } (\lambda i. A (f i)) I$
<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* $M' X' (f ' I)$
shows *indep-vars* $(M' \circ f) (\lambda k \omega. X' (f k) \omega) I$
 \langle *proof* \rangle

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

definition *covariance where*
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 [*measurable*]: $f \in \text{borel-measurable } M g \in \text{borel-measurable } M$
assumes *sq-int*: *integrable* $M (\lambda \omega. f \omega^2)$ *integrable* $M (\lambda \omega. g \omega^2)$
shows *integrable* $M (\lambda \omega. f \omega * g \omega)$
 \langle *proof* \rangle

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

lemma *sum-square-int*:
fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$
assumes $\text{finite } I$
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 $\text{integrable } M (\lambda\omega. (\sum i \in I. f i \omega)^2)$
 <proof>

lemma *var-sum-1*:
fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$
assumes $\text{finite } I$
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
 $\text{variance } (\lambda\omega. (\sum i \in I. f i \omega)) = (\sum i \in I. (\sum j \in I. \text{covariance } (f i) (f j)))$
 <proof>

lemma *covar-self-eq*:
fixes $f :: 'a \Rightarrow \text{real}$
shows $\text{covariance } f f = \text{variance } f$
 <proof>

lemma *covar-indep-eq-zero*:
fixes $f g :: 'a \Rightarrow \text{real}$
assumes $\text{integrable } M f$
assumes $\text{integrable } M g$
assumes $\text{indep-var borel } f \text{ borel } g$
shows $\text{covariance } f g = 0$
 <proof>

lemma *var-sum-2*:
fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$
assumes $\text{finite } I$
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 $\text{variance } (\lambda\omega. (\sum i \in I. f i \omega)) =$
 $(\sum i \in I. \text{variance } (f i)) + (\sum i \in I. \sum j \in I - \{i\}. \text{covariance } (f i) (f j))$
 <proof>

lemma *var-sum-pairwise-indep*:
fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$
assumes $\text{finite } I$
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 $\text{variance } (\lambda\omega. (\sum i \in I. f i \omega)) = (\sum i \in I. \text{variance } (f i))$
 <proof>

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

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

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

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

lemma *var-sum-all-indep*:

fixes $f :: 'b \Rightarrow 'a \Rightarrow \text{real}$
assumes $\text{finite } I$
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 $\text{indep-vars } (\lambda -. \text{borel}) f I$
shows $\text{variance } (\lambda\omega. (\sum i \in I. f i \omega)) = (\sum i \in I. \text{variance } (f i))$
 <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

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

lemma *pmf-prod-pmf*:

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

<proof>

lemma *PiE-default-undefined-eq*: $PiE\text{-dflt } I \text{ undefined } M = PiE \ I \ M$
<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-space.indep-vars* (*measure-pmf* (*Pi-pmf I d p*)) *M'* ($\lambda i \omega. X' i \omega$) *J*
 ⟨*proof*⟩

lemma

fixes *f* :: 'b ⇒ ('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*
 ⟨*proof*⟩

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 ⇒ 'b ⇒ ('c :: {*second-countable-topology,banach,real-normed-field*})
assumes *finite I*
assumes $\bigwedge i. i \in I \implies \text{integrable } (\text{measure-pmf } (M i)) (f i)$
shows *integral^L* (*Pi-pmf I d M*) ($\lambda x. (\prod i \in I. f i (x i))$) = $(\prod i \in I. \text{integral^L$
 (*M i*) (*f i*))
 ⟨*proof*⟩

lemma *variance-prod-pmf-slice*:

fixes *f* :: 'a ⇒ *real*
assumes $i \in I$ *finite I*
assumes *integrable* (*measure-pmf* (*M i*)) ($\lambda \omega. f \omega \hat{2}$)
shows *prob-space.variance* (*Pi-pmf I d M*) ($\lambda \omega. f (\omega i)$) = *prob-space.variance*
 (*M i*) *f*
 ⟨*proof*⟩

lemma *Pi-pmf-bind-return*:

assumes *finite I*
shows *Pi-pmf I d* ($\lambda i. M i \gg (\lambda x. \text{return-pmf } (f i x))$) = *Pi-pmf I d' M* \gg
 ($\lambda x. \text{return-pmf } (\lambda i. \text{if } i \in I \text{ then } f i (x i) \text{ else } d)$)
 ⟨*proof*⟩

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 [?]. 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-state* = *nat* × *nat* × *nat* × *nat* × (*nat* ⇒ *nat list*) × (*nat* ⇒ *float set*)

definition *hash* **where** *hash p* = *ring.hash (mod-ring p)*

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

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

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

fun *f0-space-usage* :: (*nat* × *rat* × *rat*) ⇒ *real* **where**
f0-space-usage (*n*, ε , δ) = (
 let *s* = *nat* $\lceil -18 * \ln (\text{real-of-rat } \varepsilon) \rceil$ in
 let *r* = *nat* $(4 * \lceil \log 2 (1 / \text{real-of-rat } \delta) \rceil + 23)$ in
 let *t* = *nat* $\lceil 80 / (\text{real-of-rat } \delta)^2 \rceil$ in
 6 +
 2 * *log 2* (*real s* + 1) +
 2 * *log 2* (*real t* + 1) +
 2 * *log 2* (*real n* + 21) +
 2 * *log 2* (*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 =
 $N_e \times_e (\lambda s.$
 $N_e \times_e ($
 $N_e \times_e (\lambda p.$
 $N_e \times_e ($
 $([0..<s] \rightarrow_e (P_e\ p\ 2)) \times_e$
 $([0..<s] \rightarrow_e (S_e\ F_e))))))$

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

assumes δ -range: $\delta \in \{0 < .. < 1\}$

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

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

begin

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

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

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

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** $r = nat\ (4 * \lceil \log\ 2\ (1 / real\ of\ rat\ \delta) \rceil + 23)$

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

$\langle proof \rangle$ **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*

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

private definition *sketch-rv* **where**

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

private definition *estimate*

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

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

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

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

private lemma *f0-alg-sketch*:

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

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

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

fixes $x :: nat$

fixes $q :: real$

assumes $q \geq 0$

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

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

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

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

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

assumes $c \geq 1$

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

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

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

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

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

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

private lemma

assumes $a \geq 0$

assumes $a \leq int \ p$

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

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

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

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

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

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

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

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

$as)|\} \leq 1/3$

$\langle proof \rangle$ **lemma** *median-bounds*:

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

$\langle proof \rangle$

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

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

$\langle proof \rangle$

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

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

$\langle proof \rangle$

theorem *f0-exact-space-usage*:

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

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

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

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

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

$\langle proof \rangle$

theorem *f0-asymptotic-space-complexity*:

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

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

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

$\langle proof \rangle$

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 = (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  $\delta$   $\varepsilon$  n =
    do {
      let s1 = nat ⌈6 /  $\delta^2$ ⌉;
      let s2 = nat ⌈-(18 * ln (real-of-rat  $\varepsilon$ ))⌉;
      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,  $\varepsilon$ ,  $\delta$ ) = (
    let s1 = nat ⌈6 /  $\delta^2$ ⌉ in
    let s2 = nat ⌈-(18 * ln (real-of-rat  $\varepsilon$ ))⌉ in
    3 +
    2 * log 2 (s1 + 1) +
    2 * log 2 (s2 + 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 =
 $N_e \bowtie_e (\lambda s_1.$
 $N_e \bowtie_e (\lambda s_2.$
 $N_e \bowtie_e (\lambda p.$
 $(\text{List.product } [0..<s_1] [0..<s_2] \rightarrow_e P_e p 4) \times_e$
 $(\text{List.product } [0..<s_1] [0..<s_2] \rightarrow_e I_e))))$

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

context

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

fixes $n :: \text{nat}$

fixes $as :: \text{nat list}$

fixes *result*

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

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

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

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*: $s_1 > 0$

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

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

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

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

$\langle \text{proof} \rangle$

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

$\langle \text{proof} \rangle$

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

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

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

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

$\langle \text{proof} \rangle$

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

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

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

lemma *mean-rv-alg-sketch*:

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

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

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

assumes $k < p$

shows

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

lemma

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

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

$\langle proof \rangle$

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

$\langle proof \rangle$

interpretation Ω : $prob\text{-space } \Omega_p$

$\langle proof \rangle$

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

$p)^2 - 1)$
 $\langle proof \rangle$

lemma *sketch-rv-var*:

assumes $i_2 < s_2$

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

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

lemma *mean-rv-exp*:

assumes $i < s_2$

shows $\Omega.expectation (\lambda\omega. mean-rv \omega i) = real-of-rat (F 2 as)$

$\langle proof \rangle$

lemma *mean-rv-var*:

assumes $i < s_2$

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

$\langle proof \rangle$

lemma *mean-rv-bounds*:

assumes $i < s_2$

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

lemma *f2-alg-correct'*:

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

$\langle proof \rangle$

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

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

$\langle proof \rangle$

end

Main results of this section:

theorem *f2-alg-correct*:

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

assumes $\delta > 0$

assumes $set as \subseteq \{..$

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

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

$\langle proof \rangle$

theorem *f2-exact-space-usage*:

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

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

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 $\leftarrow prod\text{-}pmf\ (\{0..<s_1\} \times \{0..<s_2\})\ (\lambda\ i.\ bernoulli\text{-}pmf\ (1 / (real\ m + 1)))$;
 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-state* \Rightarrow *rat pmf* **where**

fk-result (*s*₁, *s*₂, *k*, *m*, *r*) =
return-pmf (*median* *s*₂ ($\lambda i_2 \in \{0..<s_2\}$).
 $(\sum_{i_1 \in \{0..<s_1\}} \text{rat-of-nat } (\text{let } t = \text{snd } (r \ (i_1, i_2)) + 1 \text{ in } m * (t \hat{\sim} k - (t - 1) \hat{\sim} k))) / (\text{rat-of-nat } s_1))$)
)

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*, ε , δ) = (
let *s*₁ = *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*₂ = *nat* $\lceil -(18 * \ln (\text{real-of-rat } \varepsilon)) \rceil$ *in*
 4 +
 2 * *log 2* (*s*₁ + 1) +
 2 * *log 2* (*s*₂ + 1) +
 2 * *log 2* (*real k* + 1) +
 2 * *log 2* (*real m* + 1) +
*s*₁ * *s*₂ * (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*_e \bowtie_e (λs_1 .
*N*_e \bowtie_e (λs_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* \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* ($\lambda x s. s \gg= \text{fk-update-2 } x$) *as* (*return-pmf* (0,0,0)) =

pmf-of-set $\{..<\text{length } as\} \gg= (\lambda k. \text{return-pmf } (\text{length } as, \text{sketch } as \ k))$

<proof>

context

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

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

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 $\text{prob-space } \Omega_1$

<proof>

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

<proof>

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

<proof>

lemma

assumes $as \neq []$

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

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

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

<proof>

lemma

assumes $as \neq []$

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

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

<proof>

lemma $\text{sketch-distr}:$

assumes $as \neq []$

shows $\text{pmf-of-set } \{.. < \text{length } as\} \gg (\lambda k. \text{return-pmf } (\text{sketch } as \ k)) = \text{pmf-of-set } M_1$

<proof>

lemma *fk-update-distr*:

fold ($\lambda x s. s \gg= \text{fk-update } x$) *as* (*fk-init* $k \delta \varepsilon n$) =
prod-pmf ($\{0..<s_1\} \times \{0..<s_2\}$) ($\lambda-. \text{fold } (\lambda x s. s \gg= \text{fk-update-2 } x)$) *as* (*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 $as \neq []$
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*)
<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 $as \neq []$
shows $\text{expectation } \text{result} = \text{real-of-rat } (F \ k \ as)$

<proof>

lemma *result-var-1*:

assumes $as \neq []$

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

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

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

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

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

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

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

<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** - $\in O[?F](?rhs)$)

<proof>

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:

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

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

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 $[\cdot]_r$ 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

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.

$\{0, 1\}$ -valued random variables, with expectation $\frac{u}{p}$ and variance $\frac{u}{p} - \frac{u^2}{p^2}$.

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 $u = \lceil \frac{tp}{(1+\delta')F_0} \rceil$ 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:

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

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 $|R^* - R| \leq$

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 .

$\frac{\delta}{4}F_0$. 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:

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

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 $[h(\cdot)]_r$. 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.