Distributed Distinct Elements

Emin Karayel

May 6, 2024

Abstract

This entry formalizes a randomized cardinality estimation data structure with asymptotically optimal space usage. It is inspired by the streaming algorithm presented by Blasiok [3] in 2018. His work closed the gap between the best-known lower bound and upper bound after a long line of research started by Flajolet and Martin [4] in 1984 and was to first to apply expander graphs (in addition to hash families) to the problem. The formalized algorithm has two improvements compared to the algorithm by Blasiok. It supports operation in parallel mode, and it relies on a simpler pseudo-random construction avoiding the use of code based extractors.

Contents

1	Introduction	2
2	Preliminary Results	2
3	Blind	4
4	Balls and Bins	5
5	Tail Bounds for Expander Walks	11
6	Inner Algorithm	12
7	Accuracy without cutoff	18
8	Cutoff Level	21
9	Accuracy with cutoff	22
10	Outer Algorithm	23

1 Introduction

The algorithm is described as functional data structures, given a seed which needs to be choosen uniformly from a initial segment of the natural numbers and globally, there are three functions:

- single given the seed and an element from the universe computes a sketch for that singleton set
- merge computes a sketch based on two input sketches and returns a sketch representing the union set
- estimate computes an estimate for the cardinality of the set represented by a sketch

The main point is that a sketch requires $\mathcal{O}(\delta^{-2}\ln(\varepsilon^{-1}) + \ln n)$ space where *n* is the universe size, δ is the desired relative accuracy and ε is the desired failure probability. Note that it is easy to see that an exact solution would necessarily require $\mathcal{O}(n)$ bits.

The algorithm is split into two parts an inner algorithm, described in Section 6, which itself is already a full cardinality estimation algorithm, however its space usage is below optimal. The outer algorithm is introduced in Section 10, which runs multiple copies of the inner algorithm with carefully chosen inner parameters.

As mentioned in the abstract the algorithm is inspired by the solution to the streaming version of the problem by Błasiok [3] in 2020. His work builds on a long line of reasarch starting in 1985 [4, 1, 2, 7, 11, 5].

In an earlier AFP entry [9] I have formalized an earlier cardinality estimation algorithm based on the work by Bar-Yossef et al. [2] in 2002. Since then I have addressed the existence of finite fields for higher prime powers and expander graphs [8, 10]. Building on these results, the formalization of this more advanced solution presented here became possible.

The solution described here improves on the algorithms described by Blasiok in two ways (without comprising its optimal space usage). It can be used in a parallel mode of operation. Moreover the pseudo-random construction used is simpler than the solution described by Blasiok — who uses an extractor based on Parvaresh-Vardy codes [6] to sample random walks in an expander graph, which are then sub-sampled and then the walks are used to sample seeds for hash functions. In the solution presented here neither the sub-sampling step nor the extractor is needed, instead a two-stage expander construction is used, this means that the nodes of the first expander correspond to the walks in a second expander graph. The latters nodes correspond to seeds of hash functions (as in Blasiok's solution).

The modification needed to support a parallel mode of operation is a change in the failure strategy of the solution presented in Kane et al., which is the event when the data in the sketch reequires too much space. The main issue is that in the parallel case the number of states the algorithm might reach is not bounded by the universe size and thus an estimate they make for the probability of the failure event does not transfer to the parallel case. To solve that the algorithm in this work is more conservative. Instead of failing out-right it instead increases a cutoff threshold. For which it is then possible to show an upper estimate independent of the number of reached states.

2 Preliminary Results

This section contains various short preliminary results used in the sections below.

theory Distributed-Distinct-Elements-Preliminary imports Frequency-Moments.Frequency-Moments-Preliminary-Results Universal-Hash-Families.Universal-Hash-Families-More-Product-PMF Median-Method.Median Expander-Graphs.Extra-Congruence-Method Expander-Graphs.Constructive-Chernoff-Bound Frequency-Moments.Landau-Ext Stirling-Formula.Stirling-Formula begin unbundle intro-cong-syntax

lemma pmf-rev-mono: **assumes** $\bigwedge x. \ x \in set-pmf \ p \Longrightarrow x \notin Q \Longrightarrow x \notin P$ **shows** measure $p \ P \leq measure \ p \ Q$ $\langle proof \rangle$

lemma pmf-exp-mono: **fixes** $f g :: 'a \Rightarrow real$ **assumes** integrable (measure-pmf p) f integrable (measure-pmf p) g **assumes** $\bigwedge x. x \in set-pmf p \implies f x \leq g x$ **shows** integral^L (measure-pmf p) $f \leq integral^L$ (measure-pmf p) g $\langle proof \rangle$

lemma pmf-markov: **assumes** integrable (measure-pmf p) f c > 0 **assumes** $\bigwedge x. x \in set-pmf p \implies f x \ge 0$ **shows** measure $p \{\omega. f \omega \ge c\} \le (\int \omega. f \omega \partial p) / c$ (is $?L \le ?R$) $\langle proof \rangle$

```
lemma pair-pmf-prob-left:
measure-pmf.prob (pair-pmf A B) {\omega. P (fst \omega)} = measure-pmf.prob A {\omega. P \omega} (is ?L = ?R) 
(proof)
```

lemma pmf-exp-of-fin-function: **assumes** finite $A \ g$ ' set-pmf $p \subseteq A$ **shows** $(\int \omega. f \ (g \ \omega) \ \partial p) = (\sum y \in A. f \ y * measure \ p \ \{\omega. g \ \omega = y\})$ (is ?L = ?R) $\langle proof \rangle$

Cardinality rules for distinct/ordered pairs of a set without the finiteness constraint - to use in simplification:

lemma card-distinct-pairs: card $\{x \in B \times B. \text{ fst } x \neq \text{ snd } x\} = \text{card } B^2 - \text{card } B \text{ (is card } ?L = ?R)$ $\langle \text{proof} \rangle$ **include** intro-cong-syntax $\langle \text{proof} \rangle$ **lemma** card-ordered-pairs': **fixes** M :: ('a :: linorder) set

fixes M :: ('a ::inoraer') set **shows** card $\{(x,y) \in M \times M. \ x < y\} = card M * (card M - 1) / 2$ $\langle proof \rangle$

The following are versions of the mean value theorem, where the interval endpoints may be reversed.

lemma MVT-symmetric: **assumes** $\bigwedge x$. [[min $a \ b \le x$; $x \le max \ a \ b$]] $\Longrightarrow DERIV f x :> f' x$ **shows** $\exists z :: real. min \ a \ b \le z \land z \le max \ a \ b \land (f \ b - f \ a = (b - a) * f' z)$ $\langle proof \rangle$

```
lemma MVT-interval:

fixes I :: real set

assumes interval I a \in I \ b \in I

assumes \bigwedge x. \ x \in I \implies DERIV \ f \ x :> f' \ x

shows \exists z. \ z \in I \land (f \ b - f \ a = (b - a) * f' \ z)

\langle proof \rangle
```

Ln is monotone on the positive numbers and thus commutes with min and max:

lemma ln-min-swap: $x > (0::real) \Longrightarrow (y > 0) \Longrightarrow ln (min x y) = min (ln x) (ln y)$ $\langle proof \rangle$

lemma ln-max-swap: $x > (0::real) \Longrightarrow (y > 0) \Longrightarrow ln (max x y) = max (ln x) (ln y)$ $\langle proof \rangle$

Loose lower bounds for the factorial fuction:.

lemma fact-lower-bound: $sqrt(2*pi*n)*(n/exp(1)) \ n \leq fact \ n \ (is \ ?L \leq \ ?R) \ (proof)$

lemma fact-lower-bound-1: assumes n > 0shows $(n/exp \ 1) \ n \le fact \ n$ (is $?L \le ?R$) $\langle proof \rangle$

Rules to handle O-notation with multiple variables, where some filters may be towards zero:

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

 $\begin{array}{l} \textbf{lemma bigo-prod-1:} \\ \textbf{assumes } (\lambda x. \ f \ x) \in O[F](\lambda x. \ g \ x) \ G \neq bot \\ \textbf{shows } (\lambda x. \ f \ (fst \ x)) \in O[F \times_F \ G](\lambda x. \ g \ (fst \ x)) \\ \langle proof \rangle \end{array}$

lemma bigo-prod-2: assumes $(\lambda x. f x) \in O[G](\lambda x. g x) F \neq bot$ shows $(\lambda x. f (snd x)) \in O[F \times_F G](\lambda x. g (snd x))$ $\langle proof \rangle$

```
lemma eventually-inv:

fixes P :: real \Rightarrow bool

assumes eventually (\lambda x. P (1/x)) at-top

shows eventually (\lambda x. P x) (at-right 0)

(proof)
```

lemma bigo-inv: **fixes** $f g :: real \Rightarrow real$ **assumes** $(\lambda x. f (1/x)) \in O(\lambda x. g (1/x))$ **shows** $f \in O[at\text{-right } 0](g)$ $\langle proof \rangle$

unbundle no-intro-cong-syntax

3 Blind

Blind section added to preserve section numbers end

4 Balls and Bins

The balls and bins model describes the probability space of throwing r balls into b bins. This section derives the expected number of bins hit by at least one ball, as well as the variance in the case that each ball is thrown independently. Further, using an approximation argument it is then possible to derive bounds for the same measures in the case when the balls are being thrown only k-wise independently. The proofs follow the reasoning described in [7, §A.1] but improve on the constants, as well as constraints.

```
theory Distributed-Distinct-Elements-Balls-and-Bins
  imports
   Distributed-Distinct-Elements-Preliminary
   Discrete-Summation.Factorials
   HOL-Combinatorics.Stirling
   HOL-Computational-Algebra. Polynomial
   HOL-Decision-Procs. Approximation
begin
hide-fact Henstock-Kurzweil-Integration.integral-sum
hide-fact Henstock-Kurzweil-Integration.integral-mult-right
hide-fact Henstock-Kurzweil-Integration.integral-nonneg
hide-fact Henstock-Kurzweil-Integration.integral-cong
unbundle intro-cong-syntax
lemma sum-power-distrib:
  fixes f :: 'a \Rightarrow real
  assumes finite R
  shows (\sum i \in R. f i) \cap s = (\sum xs \mid set xs \subseteq R \land length xs = s. (\prod x \leftarrow xs. f x))
\langle proof \rangle
lemma sum-telescope-eq:
  fixes f :: nat \Rightarrow 'a :: \{comm-ring-1\}
  shows (\sum k \in \{Suc \ m.n\}, fk - f(k - 1)) = of-bool(m \le n) * (fn - fm)
  \langle proof \rangle
An improved version of diff-power-eq-sum.
lemma power-diff-sum:
  fixes a b :: 'a :: \{comm-ring-1, power\}
  shows a\hat{k} - b\hat{k} = (a-b) * (\sum i = 0..< k. a \hat{i} * b (k-1-i))
\langle proof \rangle
lemma power-diff-est:
  assumes (a :: real) \ge b
  assumes b \ge 0
  shows a^k - b^k \le (a-b) * k * a^{(k-1)}
\langle proof \rangle
lemma power-diff-est-2:
  assumes (a :: real) \ge b
  assumes b \ge 0
  shows a^k - b^k > (a-b) * k * b^{(k-1)}
\langle proof \rangle
lemma of-bool-prod:
  assumes finite R
  shows (\prod j \in R. of-bool(f j)) = (of-bool(\forall j \in R. f j) :: real)
  \langle proof \rangle
```

Additional results about falling factorials: **lemma** *ffact-nonneg*: fixes x :: realassumes $k - 1 \leq x$ shows ffact $k \ x \ge 0$ $\langle proof \rangle$ lemma ffact-pos: fixes x :: realassumes k - 1 < xshows ffact $k \ x > 0$ $\langle proof \rangle$ **lemma** *ffact-mono*: fixes x y :: realassumes $k-1 \leq x x \leq y$ **shows** flact $k \ x \leq \text{flact} \ k \ y$ $\langle proof \rangle$ **lemma** *ffact-of-nat-nonneg*: fixes $x :: 'a :: \{ comm-ring-1, linordered-nonzero-semiring \}$ assumes $x \in \mathbb{N}$ shows ffact $k \ x \ge 0$ $\langle proof \rangle$ **lemma** *ffact-suc-diff*: fixes x :: ('a :: comm-ring-1)shows flact k x - flact k (x-1) = of-nat k * flact (k-1) (x-1) (is ?L = ?R) $\langle proof \rangle$ lemma ffact-bound: ffact k (n::nat) $\leq n \hat{k}$ $\langle proof \rangle$ **lemma** *fact-moment-binomial*: fixes n :: nat and $\alpha :: real$ assumes $\alpha \in \{0...1\}$ **defines** $p \equiv binomial-pmf \ n \ \alpha$ shows $(\int \omega. \text{ flact } s \text{ (real } \omega) \partial p) = \text{flact } s \text{ (real } n) * \alpha \hat{s} \text{ (is } ?L = ?R)$ $\langle proof \rangle$ The following describes polynomials of a given maximal degree as a subset of the functions, similar to the subsets \mathbb{Z} or \mathbb{Q} as subsets of larger number classes.

definition Polynomials (\mathbb{P}) where Polynomials $k = \{f. \exists p. f = poly p \land degree p \leq k\}$

lemma Polynomials-mono: **assumes** $s \leq t$ **shows** $\mathbb{P} s \subseteq \mathbb{P} t$ $\langle proof \rangle$ **lemma** Polynomials-addI: **assumes** $f \in \mathbb{P} k \ g \in \mathbb{P} k$ **shows** $(\lambda \omega. f \ \omega + g \ \omega) \in \mathbb{P} k$ $\langle proof \rangle$

lemma Polynomials-diffI: fixes $f g :: 'a :: comm-ring \Rightarrow 'a$

assumes $f \in \mathbb{P} \ k \ g \in \mathbb{P} \ k$ shows $(\lambda x. f x - g x) \in \mathbb{P} k$ $\langle proof \rangle$ **lemma** *Polynomials-idI*: $(\lambda x. x) \in (\mathbb{P} \ 1 :: ('a::comm-ring-1 \Rightarrow 'a) \ set)$ $\langle proof \rangle$ **lemma** *Polynomials-constI*: $(\lambda x. c) \in \mathbb{P} k$ $\langle proof \rangle$ lemma Polynomials-multI: fixes $f g :: 'a :: \{comm\text{-ring}\} \Rightarrow 'a$ assumes $f \in \mathbb{P} \ s \ g \in \mathbb{P} \ t$ shows $(\lambda x. f x * g x) \in \mathbb{P}(s+t)$ $\langle proof \rangle$ **lemma** *Polynomials-composeI*: **fixes** $f g :: 'a :: \{ comm-semiring-0, semiring-no-zero-divisors \} \Rightarrow 'a$ assumes $f \in \mathbb{P} \ s \ g \in \mathbb{P} \ t$ shows $(\lambda x. f(g x)) \in \mathbb{P}(s*t)$ $\langle proof \rangle$ **lemma** *Polynomials-const-left-multI*: fixes $c :: 'a :: \{comm-ring\}$ assumes $f \in \mathbb{P} \ k$ shows $(\lambda x. \ c * f x) \in \mathbb{P} \ k$ $\langle proof \rangle$ **lemma** *Polynomials-const-right-multI*: fixes $c :: 'a :: \{comm-ring\}$ assumes $f \in \mathbb{P}$ k shows $(\lambda x. f x * c) \in \mathbb{P} k$ $\langle proof \rangle$ **lemma** *Polynomials-const-divI*: fixes $c :: 'a :: \{field\}$ assumes $f \in \mathbb{P} \ k$ shows $(\lambda x. f x / c) \in \mathbb{P} k$ $\langle proof \rangle$ **lemma** Polynomials-ffact: $(\lambda x. ffact \ s \ (x - y)) \in (\mathbb{P} \ s :: ('a :: comm-ring-1 \Rightarrow 'a) \ set)$ $\langle proof \rangle$ **lemmas** *Polynomials-intros* = Polynomials-const-divI Polynomials-composeI Polynomials-const-left-multI Polynomials-const-right-multI

Polynomials-multI

Polynomials-addI

Polynomials-diffI Polynomials-idI

Polynomials-constI

Polynomials-ffact

definition C_2 :: real where $C_2 = 7.5$

definition C_3 :: real where $C_3 = 16$

A locale fixing the sets of balls and bins

locale balls-and-bins-abs = fixes R :: 'a set and B :: 'b setassumes fin-B: finite B and B-ne: $B \neq \{\}$ assumes fin-R: finite R begin

Independent balls and bins space:

```
definition \Omega
where \Omega = prod-pmf R (\lambda-. pmf-of-set B)
```

lemma set-pmf- Ω : set-pmf $\Omega = R \rightarrow_E B$ $\langle proof \rangle$

lemma card-B-gt-0: card B > 0 $\langle proof \rangle$

```
lemma card-B-ge-1: card B \ge 1
\langle proof \rangle
```

definition $Z j \omega = real (card \{i. i \in R \land \omega i = (j::'b)\})$ definition $Y \omega = real (card (\omega ' R))$ definition $\mu = real (card B) * (1 - (1 - 1/real (card B))^card R)$

Factorial moments for the random variable describing the number of times a bin will be hit:

lemma fact-moment-balls-and-bins: **assumes** $J \subseteq B \ J \neq \{\}$ **shows** $(\int \omega. \ ffact \ s \ (\sum j \in J. \ Z \ j \ \omega) \ \partial\Omega) =$ $ffact \ s \ (real \ (card \ R)) \ * \ (real \ (card \ J) \ / \ real \ (card \ B)) \ s$ (is ?L = ?R) $\langle proof \rangle$

Expectation and variance for the number of distinct bins that are hit by at least one ball in the fully independent model. The result for the variance is improved by a factor of 4 w.r.t. the paper.

lemma

shows exp-balls-and-bins: measure-pmf.expectation $\Omega \ Y = \mu$ (is ?AL = ?AR) and var-balls-and-bins: measure-pmf.variance $\Omega \ Y \le card \ R * (real (card \ R) - 1) / card \ B$ (is ? $BL \le ?BR$)

 $\langle proof \rangle$

definition lim-balls-and-bins k p = (prob-space.k-wise-indep-vars (measure-pmf p) k (λ -. discrete) ($\lambda x \ \omega \ \omega \ x$) $R \land (\forall x. x \in R \longrightarrow map-pmf (\lambda \omega. \omega \ x) p = pmf-of-set B))$

lemma indep:

assumes lim-balls-and-bins k p shows prob-space.k-wise-indep-vars (measure-pmf p) k (λ -. discrete) ($\lambda x \ \omega \ \omega \ x$) R $\langle proof \rangle$

lemma ran: **assumes** lim-balls-and-bins $k \ p \ x \in R$ **shows** map-pmf ($\lambda \omega$. $\omega \ x$) p = pmf-of-set B $\langle proof \rangle$

lemma Z-integrable: fixes $f :: real \Rightarrow real$ assumes $lim-balls-and-bins \ k \ p$ shows integrable $p(\lambda \omega, f(Z \ i \ \omega))$ $\langle proof \rangle$ **lemma** Z-any-integrable-2: fixes $f :: real \Rightarrow real$ assumes $lim-balls-and-bins \ k \ p$ shows integrable $p (\lambda \omega. f (Z i \omega + Z j \omega))$ $\langle proof \rangle$ **lemma** *hit-count-prod-exp*: assumes $j1 \in B$ $j2 \in B$ s+t < kassumes $lim-balls-and-bins \ k \ p$ **defines** $L \equiv \{(xs, ys) \text{. set } xs \subseteq R \land set ys \subseteq R \land$ $(set \ xs \cap set \ ys = \{\} \lor j1 = j2) \land length \ xs = s \land length \ ys = t\}$ shows $(\int \omega. Z j1 \ \omega \hat{s} * Z j2 \ \omega \hat{t} \ \partial p) =$ $(\sum (xs, ys) \in L. (1/real (card B)) \cap (card (set xs \cup set ys)))$ $(\mathbf{is} ?L = ?R)$ $\langle proof \rangle$ **lemma** *hit-count-prod-pow-eq*: assumes $i \in B \ j \in B$ **assumes** lim-balls-and-bins k passumes lim-balls-and-bins k qassumes $s+t \leq k$ shows $(\int \omega. (Z \ i \ \omega) \ s \ast (Z \ j \ \omega) \ t \ \partial p) = (\int \omega. (Z \ i \ \omega) \ s \ast (Z \ j \ \omega) \ t \ \partial q)$ $\langle proof \rangle$ **lemma** *hit-count-sum-pow-eq*: assumes $i \in B \ j \in B$ **assumes** *lim-balls-and-bins* k passumes *lim-balls-and-bins* k q assumes s < kshows $(\int \omega. (Z \, i \, \omega + Z \, j \, \omega) \, \hat{s} \, \partial p) = (\int \omega. (Z \, i \, \omega + Z \, j \, \omega) \, \hat{s} \, \partial q)$ (is ?L = ?R) $\langle proof \rangle$ **lemma** *hit-count-sum-poly-eq*: assumes $i \in B \ j \in B$ assumes lim-balls-and-bins k passumes $lim-balls-and-bins \ k \ q$ assumes $f \in \mathbb{P} \ k$ shows $(\int \omega. f (Z i \omega + Z j \omega) \partial p) = (\int \omega. f (Z i \omega + Z j \omega) \partial q)$ (is ?L = ?R) $\langle proof \rangle$ **lemma** *hit-count-poly-eq*: assumes $b \in B$ assumes $lim-balls-and-bins \ k \ p$ **assumes** *lim-balls-and-bins* k qassumes $f \in \mathbb{P} \ k$ shows $(\int \omega. f (Z \ b \ \omega) \ \partial p) = (\int \omega. f (Z \ b \ \omega) \ \partial q)$ (is ?L = ?R) $\langle proof \rangle$

 ${\bf lemma}\ lim-balls-and-bins-from-ind-balls-and-bins:$

$$\begin{split} & lim-balls-and-bins \ k \ \Omega \\ & \langle proof \rangle \end{split}$$
 $\begin{aligned} & lemma \ hit-count-factorial-moments: \\ & assumes \ a:j \in B \\ & assumes \ s \leq k \\ & assumes \ lim-balls-and-bins \ k \ p \\ & shows \ (\int \omega. \ ffact \ s \ (Z \ j \ \omega) \ \partial p) = ffact \ s \ (real \ (card \ R)) \ * \ (1 \ / \ real \ (card \ B))^{\ s} \\ & (is \ ?L = \ ?R) \\ & \langle proof \rangle \end{aligned}$

lemma hit-count-factorial-moments-2: **assumes** $a:i \in B \ j \in B$ **assumes** $i \neq j \ s \leq k \ card \ R \leq card \ B$ **assumes** $lim-balls-and-bins \ k \ p$ **shows** $(\int \omega. \ ffact \ s \ (Z \ i \ \omega + Z \ j \ \omega) \ \partial p) \leq 2^s$ (is $?L \leq ?R)$ $\langle proof \rangle$

```
lemma balls-and-bins-approx-helper:

fixes x :: real

assumes x \ge 2

assumes real k \ge 5*x / \ln x

shows k \ge 2

and 2 \widehat{(k+3)} / fact k \le (1/exp x) \widehat{2}

and 2 / fact k \le 1 / (exp \ 1 * exp \ x)

\langle proof \rangle
```

Bounds on the expectation and variance in the k-wise independent case. Here the indepedence assumption is improved by a factor of two compared to the result in the paper.

lemma

assumes card $R \leq card B$ assumes $\bigwedge c. \ lim-balls-and-bins \ (k+1) \ (p \ c)$ assumes $\varepsilon \in \{0 < ...1/exp(2)\}$ assumes $k \geq 5 * \ln (card B / \varepsilon) / \ln (\ln (card B / \varepsilon))$ shows $exp-approx: |measure-pmf.expectation \ (p \ True) \ Y - measure-pmf.expectation \ (p \ False) \ Y| \leq \varepsilon * real \ (card \ R) \ (is \ ?A) \ and$ $var-approx: |measure-pmf.variance \ (p \ True) \ Y - measure-pmf.variance \ (p \ False) \ Y| \leq \varepsilon^2$ $(is \ ?B)$ $\langle proof \rangle$

lemma

assumes card $R \leq card B$ assumes lim-balls-and-bins (k+1) passumes $k \geq 7.5 * (ln (card B) + 2)$ shows exp-approx-2: |measure-pmf.expectation $p Y - \mu| \leq card R / sqrt (card B)$ (is $?AL \leq ?AR$) and var-approx-2: measure-pmf.variance $p Y \leq real (card R)^2 / card B$ (is $?BL \leq ?BR$) $\langle proof \rangle$

lemma devitation-bound: **assumes** card $R \le card B$ **assumes** lim-balls-and-bins k p **assumes** real $k \ge C_2 * ln (real (card B)) + C_3$ **shows** measure $p \{\omega. | Y \omega - \mu| > 9 * real (card R) / sqrt (real (card B)) \} \le 1 / 2^6$ (is $?L \le ?R$) $\langle proof \rangle$

 \mathbf{end}

unbundle no-intro-cong-syntax

end

5 Tail Bounds for Expander Walks

theory Distributed-Distinct-Elements-Tail-Bounds imports Distributed-Distinct-Elements-Preliminary Expander-Graphs.Pseudorandom-Objects-Expander-Walks HOL-Decision-Procs.Approximation

begin

This section introduces tail estimates for random walks in expander graphs, specific to the verification of this algorithm (in particular to two-stage expander graph sampling and obtained tail bounds for subgaussian random variables). They follow from the more fundamental results *regular-graph.kl-chernoff-property* and *regular-graph.uniform-property* which are verified in the AFP entry for expander graphs [10].

 ${\bf hide-fact} \ {\it Henstock-Kurzweil-Integration.integral-sum}$

```
unbundle intro-cong-syntax
```

lemma *x*-ln-*x*-min: assumes $x \ge (0::real)$ shows $x * \ln x \ge -exp(-1)$ $\langle proof \rangle$ theorem (in regular-graph) walk-tail-bound: assumes l > 0**assumes** $S \subseteq verts G$ defines $\mu \equiv real (card S) / card (verts G)$ assumes $\gamma < 1 \ \mu + \Lambda_a \leq \gamma$ shows measure (pmf-of-multiset (walks G l)) {w. real (card $\{i \in \{... < l\}, w \mid i \in S\}$) $\geq \gamma * l$ } $\leq exp \ (-real \ l * (\gamma * ln \ (1/(\mu + \Lambda_a)) - 2 * exp(-1))) \ (is \ ?L \leq ?R)$ $\langle proof \rangle$ **theorem** (in regular-graph) walk-tail-bound-2: assumes l > 0 $\Lambda_a \leq \Lambda \Lambda > 0$ assumes $S \subseteq verts G$ defines $\mu \equiv real (card S) / card (verts G)$ assumes $\gamma < 1 \ \mu + \Lambda \leq \gamma$ shows measure (pmf-of-multiset (walks G l)) {w. real (card $\{i \in \{..< l\}. w \mid i \in S\}$) $\geq \gamma * l$ } $\leq exp \ (-real \ l * (\gamma * ln \ (1/(\mu+\Lambda)) - 2 * exp(-1))) \ (is \ ?L \leq ?R)$ $\langle proof \rangle$ **lemma** disjI-safe: $(\neg x \Longrightarrow y) \Longrightarrow x \lor y \langle proof \rangle$

lemma walk-tail-bound: fixes T assumes l > 0 $\Lambda > 0$ assumes measure (sample-pro S) $\{w. T w\} \le \mu$ assumes $\gamma \le 1 \ \mu + \Lambda \le \gamma \ \mu \le 1$ shows measure (sample-pro ($\mathcal{E} \ l \ \Lambda S$)) $\{w. real (card \ \{i \in \{..< l\}. T (w i)\}) \ge \gamma * l\}$ $\leq exp \ (- \ real \ l * (\gamma * ln \ (1/(\mu + \Lambda)) - 2 * exp(-1))) \ (is \ ?L \leq ?R) \\ \langle proof \rangle$

definition C_1 :: real where $C_1 = exp \ 2 + exp \ 3 + (exp \ 1 - 1)$

lemma deviation-bound: fixes $f :: 'a \Rightarrow real$ assumes l > 0assumes $\Lambda \in \{0 < ...exp (-real \ l * ln (real \ l)^3)\}$ assumes $\Lambda x. x \ge 20 \implies measure (sample-pro \ S) \{v. f \ v \ge x\} \le exp (-x * ln \ x^3)$ shows measure (sample-pro $(\mathcal{E} \ l \ \Lambda \ S)) \{\omega. (\sum i < l. f (\omega \ i)) \ge C_1 * l\} \le exp (-real \ l)$ (is ?L $\le ?R)$ $\langle proof \rangle$

unbundle no-intro-cong-syntax

 \mathbf{end}

6 Inner Algorithm

This section introduces the inner algorithm (as mentioned it is already a solution to the cardinality estimation with the caveat that, if ε is too small it requires to much space. The outer algorithm in Section 10 resolved this problem.

The algorithm makes use of the balls and bins model, more precisely, the fact that the number of hit bins can be used to estimate the number of balls thrown (even if there are collusions). I.e. it assigns each universe element to a bin using a k-wise independent hash function. Then it counts the number of bins hit.

This strategy however would only work if the number of balls is roughly equal to the number of bins, to remedy that the algorithm performs an adaptive sub-sampling strategy. This works by assigning each universe element a level (using a second hash function) with a geometric distribution. The algorithm then selects a level that is appropriate based on a rough estimate obtained using the maximum level in the bins.

To save space the algorithm drops information about small levels, whenever the space usage would be too high otherwise. This level will be called the cutoff-level. This is okey as long as the cutoff level is not larger than the sub-sampling threshold. A lot of the complexity in the proof is devoted to verifying that the cutoff-level will not cross it, it works by defining a third value s_M that is both an upper bound for the cutoff level and a lower bound for the sub-sampling threshold simultaneously with high probability.

theory Distributed-Distinct-Elements-Inner-Algorithm imports

Universal-Hash-Families.Pseudorandom-Objects-Hash-Families Distributed-Distinct-Elements-Preliminary Distributed-Distinct-Elements-Balls-and-Bins Distributed-Distinct-Elements-Tail-Bounds Prefix-Free-Code-Combinators.Prefix-Free-Code-Combinators begin

unbundle intro-cong-syntax hide-const Abstract-Rewriting.restrict

definition C_4 :: real where $C_4 = 3^2 * 2^2 3$ definition C_5 :: int where $C_5 = 33$ definition C_6 :: real where $C_6 = 4$ definition C_7 :: nat where $C_7 = 2^5$

locale inner-algorithm =fixes n :: natfixes δ :: real fixes ε :: real assumes *n*-gt- θ : $n > \theta$ assumes δ -gt- θ : $\delta > \theta$ and δ -lt-1: $\delta < 1$ assumes ε -gt- θ : $\varepsilon > \theta$ and ε -lt-1: $\varepsilon < 1$ begin definition b-exp where b-exp = nat $\lceil \log 2 (C_4 / \varepsilon^2) \rceil$ definition b :: nat where $b = 2^{b-exp}$ definition l where $l = nat \left[C_6 * ln \left(\frac{2}{\delta} \right) \right]$ definition k where $k = nat \left[C_2 * ln \ b + C_3 \right]$ definition Λ :: real where $\Lambda = min (1/16) (exp (-l * ln l^3))$ **definition** ρ :: real \Rightarrow real where $\rho x = b * (1 - (1 - 1/b) powr x)$ definition ρ -inv :: real \Rightarrow real where ρ -inv $x = \ln(1-x/b) / \ln(1-1/b)$ lemma *l*-lbound: $C_6 * ln (2 / \delta) \leq l$ $\langle proof \rangle$ **lemma** k-min: $C_2 * ln (real b) + C_3 \le real k$ $\langle proof \rangle$ lemma Λ -gt- θ : $\Lambda > \theta$ $\langle proof \rangle$ lemma Λ -le-1: $\Lambda \leq 1$ $\langle proof \rangle$ lemma *l-gt-0*: l > 0 $\langle proof \rangle$ lemma *l-ubound*: $l \leq C_6 * ln(1 / \delta) + C_6 * ln 2 + 1$ $\langle proof \rangle$ lemma b-exp-ge-26: b-exp ≥ 26 $\langle proof \rangle$ lemma *b-min*: $b \ge 2^2 6$ $\langle proof \rangle$ lemma k-gt-0: k > 0 $\langle proof \rangle$ **lemma** *b*-*ne*: $\{.. < b\} \neq \{\}$ $\langle proof \rangle$ **lemma** b-lower-bound: $C_4 / \varepsilon 2 \leq real b$ $\langle proof \rangle$ definition *n*-exp where n-exp = max (nat $\lceil \log 2 n \rceil$) 1 lemma *n*-exp-gt- θ : *n*-exp > θ $\langle proof \rangle$ abbreviation Ψ_1 where $\Psi_1 \equiv \mathcal{H} \ 2 \ n \ (\mathcal{G} \ n\text{-}exp)$ abbreviation Ψ_2 where $\Psi_2 \equiv \mathcal{H} \ 2 \ n \ (\mathcal{N} \ (C_7 * b^2))$

abbreviation Ψ_3 where $\Psi_3 \equiv \mathcal{H} \ k \ (C_7 * b^2) \ (\mathcal{N} \ b)$ definition Ψ where $\Psi = \Psi_1 \times_P \Psi_2 \times_P \Psi_3$ abbreviation Ω where $\Omega \equiv \mathcal{E} \ l \ \Lambda \ \Psi$ **type-synonym** state = $(nat \Rightarrow nat \Rightarrow int) \times (nat)$ fun *is-too-large* :: $(nat \Rightarrow nat \Rightarrow int) \Rightarrow bool$ where *is-too-large* $B = ((\sum (i,j) \in \{..< l\} \times \{..< b\}, \lfloor log \ 2 \ (max \ (B \ i \ j) \ (-1) + 2) \rfloor) > C_5 * b * l)$ **fun** *compress-step* :: *state* \Rightarrow *state* **where** compress-step $(B,q) = (\lambda \ i \ j. \ max \ (B \ i \ j - 1) \ (-1), \ q+1)$ **function** *compress* :: *state* \Rightarrow *state* **where** compress (B,q) = (if is-too-large Bthen (compress (compress-step (B,q))) else (B,q)) $\langle proof \rangle$ **fun** compress-termination :: state \Rightarrow nat **where** compress-termination $(B,q) = (\sum (i,j) \in \{..< l\} \times \{..< b\}$. nat $(B \ i \ j + 1))$ **lemma** compress-termination: assumes is-too-large B shows compress-termination (compress-step (B,q)) < compress-termination (B,q) $\langle proof \rangle$ termination compress $\langle proof \rangle$ **fun** merge1 :: state \Rightarrow state **where** merge1 $(B1,q_1)$ $(B2, q_2) = ($ let $q = max q_1 q_2$ in $(\lambda \ i \ j. \ max (B1 \ i \ j + q_1 - q) (B2 \ i \ j + q_2 - q), q))$ **fun** *merge* :: *state* \Rightarrow *state* \Rightarrow *state* **where** merge x y = compress (merge1 x y)**type-synonym** seed = $nat \Rightarrow (nat \Rightarrow nat) \times (nat \Rightarrow nat) \times (nat \Rightarrow nat)$ **fun** single1 :: seed \Rightarrow nat \Rightarrow state where single1 $\omega x = (\lambda \ i \ j.$ let $(f,g,h) = \omega$ i in (if $h(g x) = j \land i < l$ then int (f x) else (-1), 0) **fun** single :: seed \Rightarrow nat \Rightarrow state **where** single $\omega x = compress (single1 \ \omega x)$ **fun** *estimate1* :: *state* \Rightarrow *nat* \Rightarrow *real* **where** estimate1 (B,q) i = (

 $\begin{array}{l} \text{estimater } (D,q) \ i = (\\ \text{let } s = max \ 0 \ (Max \ ((B \ i) \ `\{..<b\}) + q - \lfloor \log 2 \ b \rfloor + 9);\\ p = card \ \{ \ j. \ j \in \{..<b\} \land B \ i \ j + q \geq s \ \} \ in\\ 2 \ powr \ s \ s \ ln \ (1-p/b) \ / \ ln(1-1/b)) \end{array}$

fun estimate :: state \Rightarrow real where estimate x = median l (estimate1 x)

6.1 History Independence

fun $\tau_0 :: ((nat \Rightarrow nat) \times (nat \Rightarrow nat) \times (nat \Rightarrow nat)) \Rightarrow nat set \Rightarrow nat \Rightarrow int$ where τ_0 (f,g,h) A = Max ({ int (f a) | a . a \in A \land h (g a) = j } $\cup \{-1\}$) definition $\tau_1 :: ((nat \Rightarrow nat) \times (nat \Rightarrow nat) \times (nat \Rightarrow nat)) \Rightarrow nat set \Rightarrow nat \Rightarrow int$ where $\tau_1 \psi A q j = max (\tau_0 \psi A j - q) (-1)$ **definition** $\tau_2 :: seed \Rightarrow nat set \Rightarrow nat \Rightarrow nat \Rightarrow nat \Rightarrow int$ where $\tau_2 \ \omega \ A \ q \ i \ j = (if \ i < l \ then \ \tau_1 \ (\omega \ i) \ A \ q \ j \ else \ (-1))$ **definition** $\tau_3 :: seed \Rightarrow nat set \Rightarrow nat \Rightarrow state$ where $\tau_3 \ \omega \ A \ q = (\tau_2 \ \omega \ A \ q, \ q)$ **definition** $q :: seed \Rightarrow nat set \Rightarrow nat$ where $q \ \omega \ A = (LEAST \ q \ . \neg (is-too-large \ (\tau_2 \ \omega \ A \ q)))$ **definition** $\tau :: seed \Rightarrow nat set \Rightarrow state$ where $\tau \ \omega \ A = \tau_3 \ \omega \ A \ (q \ \omega \ A)$ **lemma** τ_2 -step: $\tau_2 \omega A (x+y) = (\lambda i j. max (\tau_2 \omega A x i j - y) (-1))$ $\langle proof \rangle$ lemma τ_3 -step: compress-step ($\tau_3 \ \omega \ A \ x$) = $\tau_3 \ \omega \ A \ (x+1)$ $\langle proof \rangle$ lemma Ψ_1 : *is-prime-power* (*pro-size* (\mathcal{G} *n-exp*)) $\langle proof \rangle$ lemma Ψ_2 : is-prime-power (pro-size ($\mathcal{N}(C_7 * b^2)$)) $\langle proof \rangle$ lemma Ψ_3 : is-prime-power (pro-size (\mathcal{N} b)) $\langle proof \rangle$ lemma sample-pro- Ψ : sample-pro Ψ = pair-pmf (sample-pro Ψ_1) (pair-pmf (sample-pro Ψ_2) (sample-pro Ψ_3)) $\langle proof \rangle$ **lemma** sample-set- Ψ : pro-set Ψ = pro-set Ψ ₁ × pro-set Ψ ₂ × pro-set Ψ ₃ $\langle proof \rangle$ lemma *f*-range: assumes $(f,g,h) \in pro\text{-}set \Psi$ shows $f x \leq n$ -exp $\langle proof \rangle$ **lemma** *g*-range-1: assumes $g \in pro\text{-set } \Psi_2$ shows $g x < C_7 * b^2$ $\langle proof \rangle$ **lemma** *h*-range-1: assumes $h \in pro\text{-set } \Psi_3$ shows h x < b $\langle proof \rangle$ lemma g-range:

assumes $(f,g,h) \in pro\text{-set } \Psi$ shows $g x < C_7 * b^2$ $\langle proof \rangle$ **lemma** *h*-range: assumes $(f,g,h) \in pro\text{-set } \Psi$ shows h x < b $\langle proof \rangle$ lemma *fin-f*: assumes $(f,g,h) \in pro\text{-set } \Psi$ shows finite { int $(f a) \mid a. P a$ } (is finite ?M) $\langle proof \rangle$ **lemma** Max-int-range: $x \leq (y::int) \Longrightarrow Max \{x..y\} = y$ $\langle proof \rangle$ lemma $\Omega: l > \theta \Lambda > \theta \langle proof \rangle$ lemma ω -range: assumes $\omega \in pro\text{-set } \Omega$ shows $\omega \ i \in pro\text{-set } \Psi$ $\langle proof \rangle$ **lemma** *max-q-1*: assumes $\omega \in pro\text{-set } \Omega$ shows $\tau_2 \omega A (nat \lceil log \ 2 \ n \rceil + 2) i j = (-1)$ $\langle proof \rangle$ lemma *max-q-2*: assumes $\omega \in pro\text{-set } \Omega$ **shows** \neg (*is-too-large* ($\tau_2 \omega A (nat \lceil log \ 2 \ n \rceil + 2)))$ $\langle proof \rangle$ lemma max-s-3: assumes $\omega \in pro\text{-set } \Omega$ shows $q \ \omega \ A \le (nat \lceil log \ 2 \ n \rceil + 2)$ $\langle proof \rangle$ **lemma** max-mono: $x \leq (y::'a::linorder) \implies max \ x \ z \leq max \ y \ z$ $\langle proof \rangle$ **lemma** max-mono-2: $y \leq (z::'a::linorder) \implies max \ x \ y \leq max \ x \ z$ $\langle proof \rangle$ lemma τ_0 -mono: assumes $\psi \in pro\text{-set } \Psi$ **assumes** $A \subseteq B$ shows $\tau_0 \ \psi \ A \ j \leq \tau_0 \ \psi \ B \ j$ $\langle proof \rangle$ lemma τ_2 -mono: assumes $\omega \in pro\text{-set } \Omega$ assumes $A \subseteq B$ shows $\tau_2 \ \omega \ A \ x \ i \ j \leq \tau_2 \ \omega \ B \ x \ i \ j$ $\langle proof \rangle$ **lemma** *is-too-large-antimono*:

assumes $\omega \in pro\text{-set } \Omega$ **assumes** $A \subseteq B$ assumes is-too-large ($\tau_2 \ \omega \ A \ x$) shows is-too-large $(\tau_2 \ \omega \ B \ x)$ $\langle proof \rangle$ **lemma** *q*-compact: assumes $\omega \in pro\text{-set } \Omega$ **shows** \neg (*is-too-large* ($\tau_2 \ \omega \ A \ (q \ \omega \ A))$) $\langle proof \rangle$ **lemma** *q*-mono: assumes $\omega \in pro\text{-}set \ \Omega$ assumes $A \subseteq B$ shows $q \ \omega \ A \leq q \ \omega \ B$ $\langle proof \rangle$ **lemma** *lt-s-too-large*: $x < q \omega A \Longrightarrow is$ -too-large $(\tau_2 \omega A x)$ $\langle proof \rangle$ **lemma** compress-result-1: assumes $\omega \in pro\text{-set } \Omega$ shows compress $(\tau_3 \ \omega \ A \ (q \ \omega \ A - i)) = \tau \ \omega \ A$ $\langle proof \rangle$ lemma compress-result: assumes $\omega \in pro\text{-set } \Omega$ assumes $x \leq q \ \omega \ A$ shows compress $(\tau_3 \ \omega \ A \ x) = \tau \ \omega \ A$ $\langle proof \rangle$ lemma τ_0 -merge: assumes $(f,g,h) \in pro\text{-set } \Psi$ shows τ_0 (f,g,h) $(A \cup B)$ j = max $(\tau_0 (f,g,h) A j)$ $(\tau_0 (f,g,h) B j)$ (is ?L = ?R) $\langle proof \rangle$ lemma τ_2 -merge: assumes $\omega \in pro\text{-set } \Omega$ shows $\tau_2 \omega (A \cup B) x i j = max (\tau_2 \omega A x i j) (\tau_2 \omega B x i j)$ $\langle proof \rangle$ **lemma** *merge1-result*: assumes $\omega \in pro\text{-set } \Omega$ shows merge1 ($\tau \ \omega \ A$) ($\tau \ \omega \ B$) = $\tau_3 \ \omega \ (A \cup B) \ (max \ (q \ \omega \ A) \ (q \ \omega \ B))$ $\langle proof \rangle$ **lemma** *merge-result*: assumes $\omega \in pro\text{-set } \Omega$ shows merge $(\tau \ \omega \ A) \ (\tau \ \omega \ B) = \tau \ \omega \ (A \cup B)$ (is ?L = ?R) $\langle proof \rangle$ **lemma** single1-result: single1 $\omega x = \tau_3 \omega \{x\} 0$ $\langle proof \rangle$ **lemma** *single-result*: assumes $\omega \in pro\text{-set } \Omega$ shows single $\omega x = \tau \omega \{x\}$ (is ?L = ?R) $\langle proof \rangle$

6.2 Encoding states of the inner algorithm

definition *is-state-table* :: $(nat \times nat \Rightarrow int) \Rightarrow bool$ where *is-state-table* $g = (range \ g \subseteq \{-1..\} \land g \ (-(\{..< l\} \times \{..< b\})) \subseteq \{-1\})$

Encoding for state table values:

definition V_e :: int encoding where $V_e x = (if x \ge -1 then N_e (nat (x+1)) else None)$

Encoding for state table:

 $\begin{array}{l} \textbf{definition } T_e' :: (nat \times nat \Rightarrow int) \ encoding \ \textbf{where} \\ T_e' \ g = (\\ if \ is-state-table \ g \\ then \ (List.product \ [0..< l] \ [0..< b] \rightarrow_e \ V_e) \ (restrict \ g \ (\{..< l\} \times \{..< b\})) \\ else \ None) \end{array}$

definition $T_e :: (nat \Rightarrow nat \Rightarrow int)$ encoding where $T_e f = T_e'$ (case-prod f)

definition encode-state :: state encoding where encode-state = $T_e \times_e Nb_e$ (nat $\lceil \log 2 n \rceil + 3$)

lemma inj-on-restrict: **assumes** $B \subseteq \{f. f ` (-A) \subseteq \{c\}\}$ **shows** inj-on ($\lambda x.$ restrict x A) B $\langle proof \rangle$

lemma encode-state: is-encoding encode-state $\langle proof \rangle$

lemma state-bit-count: assumes $\omega \in pro\text{-set }\Omega$ shows bit-count (encode-state $(\tau \ \omega \ A)$) $\leq 2^36 * (\ln(1/\delta)+1)/\varepsilon^2 + \log 2 \ (\log 2 \ n + 3)$ (is $?L \leq ?R$) $\langle proof \rangle$

lemma random-bit-count: pro-size $\Omega \leq 2$ powr (4 * log 2 n + 48 * (log 2 (1 / ε) + 16)² + (55 + 60 * ln (1 / δ))³) (is ?L \leq ?R) (proof)

end

unbundle no-intro-cong-syntax

end

7 Accuracy without cutoff

This section verifies that each of the l estimate have the required accuracy with high probability assuming that there was no cut-off, i.e., that s = 0. Section 9 will then show that this remains true as long as the cut-off is below t f the subsampling threshold.

theory Distributed-Distinct-Elements-Accuracy-Without-Cutoff imports Concentration-Inequalities.Bienaymes-Identity Distributed-Distinct-Elements-Inner-Algorithm Distributed-Distinct-Elements-Balls-and-Bins begin

hide-fact (open) Discrete.log-mono **no-notation** Polynomials.var (X_1) locale inner-algorithm-fix-A = inner-algorithm +fixes Aassumes A-range: $A \subseteq \{..< n\}$ assumes A-nonempty: $\{\} \neq A$ begin definition X :: nat where X = card Adefinition q-max where q-max = nat ($\lceil \log 2 X \rceil - b$ -exp) **definition** $t :: (nat \Rightarrow nat) \Rightarrow int$ where t f = int (Max (f' A)) - b - exp + 9**definition** $s :: (nat \Rightarrow nat) \Rightarrow nat$ where s f = nat (t f)**definition** $R :: (nat \Rightarrow nat) \Rightarrow nat set$ where $R f = \{a. a \in A \land f a \ge s f\}$ **definition** $r :: nat \Rightarrow (nat \Rightarrow nat) \Rightarrow nat$ where $r x f = card \{a. a \in A \land f a \ge x\}$ **definition** p where $p = (\lambda(f,g,h). card \{j \in \{.., <b\}, \tau_1 (f,g,h) \land 0 j \ge s f\})$ definition Y where $Y = (\lambda(f,g,h), 2 \uparrow s f * \varrho \text{-inv} (p (f,g,h)))$ lemma fin-A: finite A $\langle proof \rangle$ lemma X-le-n: $X \leq n$ $\langle proof \rangle$ lemma X-ge-1: $X \ge 1$ $\langle proof \rangle$ **lemma** of-bool-square: $(of-bool x)^2 = ((of-bool x)::real)$ $\langle proof \rangle$ **lemma** r-eq: $r x f = (\sum a \in A.(of-bool(x \le f a) :: real))$ $\langle proof \rangle$ lemma shows *r-exp*: $(\int \omega \cdot real (r \cdot x \cdot \omega) \partial \Psi_1) = real X * (of-bool (x \le max (nat \lceil \log 2 \cdot n \rceil) 1) / 2^x)$ and *r*-var: measure-pmf.variance Ψ_1 ($\lambda \omega$. real ($r \ x \ \omega$)) $\leq (\int \omega$. real ($r \ x \ \omega$) $\partial \ \Psi_1$) $\langle proof \rangle$ definition E_1 where $E_1 = (\lambda(f,g,h), 2 \text{ powr } (-t f) * X \in \{b/2 \widehat{1}6..b/2\})$ lemma *t-low*: measure Ψ_1 {f. of-int (t f) < log 2 (real X) + 1 - b-exp} $\leq 1/2^{\gamma}$ (is $?L \leq ?R$) $\langle proof \rangle$

lemma t-high: measure Ψ_1 {f. of-int (t f) > log 2 (real X) + 16 - b-exp} $\leq 1/2^{\gamma}$ (is $?L \leq ?R$) $\langle proof \rangle$

lemma e-1: measure $\Psi \{\psi, \neg E_1 \psi\} \leq 1/2\hat{} \delta \langle proof \rangle$

definition E_2 where $E_2 = (\lambda(f,g,h). |card (R f) - X / 2^{(s f)}| \le \varepsilon/3 * X / 2^{(s f)})$

lemma e-2: measure Ψ { ψ . $E_1 \ \psi \land \neg E_2 \ \psi$ } $\leq 1/2^6$ (is $?L \leq ?R$) $\langle proof \rangle$

definition E_3 where $E_3 = (\lambda(f,g,h). inj on g (R f))$

lemma *R*-bound: fixes f g hassumes $E_1 (f,g,h)$ assumes $E_2 (f,g,h)$ shows card $(R f) \le 2/3 * b$ $\langle proof \rangle$

lemma e-3: measure Ψ { ψ . $E_1 \psi \wedge E_2 \psi \wedge \neg E_3 \psi$ } $\leq 1/2^6$ (is $?L \leq ?R$) \(\lapproof\)

definition E_4 where $E_4 = (\lambda(f,g,h). |p(f,g,h) - \varrho(card(R f))| \le \varepsilon/12 * card(R f))$

lemma e-4-h: 9 / sqrt $b \leq \varepsilon$ / 12 $\langle proof \rangle$

lemma e-4: measure Ψ { ψ . $E_1 \psi \wedge E_2 \psi \wedge E_3 \psi \wedge \neg E_4 \psi$ } $\leq 1/2^6$ (is $?L \leq ?R$) $\langle proof \rangle$

lemma ρ -inverse: ρ -inv $(\rho x) = x$ $\langle proof \rangle$

lemma rho-mono: assumes $x \le y$ shows $\varrho \ x \le \varrho \ y$ $\langle proof \rangle$

lemma rho-two-thirds: $\rho (2/3 * b) \leq 3/5 * b \langle proof \rangle$

definition ρ -inv' :: real \Rightarrow real where ρ -inv' $x = -1 / (real \ b * (1-x / real \ b) * ln (1 - 1 / real \ b))$

```
\begin{array}{ll} \textbf{lemma } \varrho \text{-}inv': \\ \textbf{fixes } x :: real \\ \textbf{assumes } x < b \\ \textbf{shows } DERIV \ \varrho \text{-}inv \ x :> \varrho \text{-}inv' \ x \\ \langle proof \rangle \end{array}
```

```
lemma accuracy-without-cutoff:
measure \Psi \{(f,g,h). | Y (f,g,h) - real X | > \varepsilon * X \lor s f < q-max\} \le 1/2^4
(is ?L \le ?R)
\langle proof \rangle
```

end

 \mathbf{end}

8 Cutoff Level

This section verifies that the cutoff will be below q-max with high probability. The result will be needed in Section 9, where it is shown that the estimates will be accurate for any cutoff below q-max.

theory Distributed-Distinct-Elements-Cutoff-Level

```
imports
Distributed-Distinct-Elements-Accuracy-Without-Cutoff
Distributed-Distinct-Elements-Tail-Bounds
```

begin

hide-const (open) Quantum.Z

unbundle intro-cong-syntax

lemma mono-real-of-int: mono real-of-int $\langle proof \rangle$

lemma Max-le-Sum: **fixes** $f :: a \Rightarrow int$ **assumes** finite A **assumes** $\bigwedge a. \ a \in A \implies f \ a \ge 0$ **shows** Max (insert 0 (f ' A)) $\le (\sum a \in A . f \ a)$ (is $?L \le ?R$) $\langle proof \rangle$

context inner-algorithm-fix-A **begin**

The following inequality is true for base e on the entire domain (x > 0). It is shown in *ln-add-one-self-le-self*. In the following it is established for base 2, where it holds for $x \ge 1$.

```
\begin{array}{l} \textbf{lemma } log-2\text{-estimate:} \\ \textbf{assumes } x \geq (1::real) \\ \textbf{shows } log \ 2 \ (1+x) \leq x \\ \langle proof \rangle \end{array}
\begin{array}{l} \textbf{lemma } cutoff\text{-}eq\text{-}7\text{:} \\ real \ X * 2 \ powr \ (-real \ q\text{-}max) \ / \ b \leq 1 \\ \langle proof \rangle \end{array}
\begin{array}{l} \textbf{lemma } cutoff\text{-}eq\text{-}6\text{:} \\ \textbf{fixes } k \\ \textbf{assumes } a \in A \\ \textbf{shows } \ (\int f. \ real\text{-}of\text{-}int \ (max \ 0 \ (int \ (f \ a) \ - \ int \ k)) \ \partial \Psi_1) \leq 2 \ powr \ (-real \ k) \ (\textbf{is } \ ?L \leq \ ?R) \\ \langle proof \rangle \end{array}
\begin{array}{l} \textbf{lemma } cutoff\text{-}eq\text{-}5\text{:} \end{array}
```

```
assumes x \ge (-1 :: real)
```

shows real-of-int $|\log 2(x+2)| \leq (real c+2) + max(x-2^{c}) 0$ (is $?L \leq ?R$) $\langle proof \rangle$

```
lemma cutoff-level:
  measure \Omega {\omega. q \ \omega \ A > q-max} \leq \delta/2 (is ?L \leq ?R)
\langle proof \rangle
```

end

unbundle no-intro-cong-syntax

end

9 Accuracy with cutoff

This section verifies that each of the l estimate have the required accuracy with high probability assuming as long as the cutoff is below q-max, generalizing the result from Section 7.

```
theory Distributed-Distinct-Elements-Accuracy
 imports
   Distributed-Distinct-Elements-Accuracy-Without-Cutoff
   Distributed-Distinct-Elements-Cutoff-Level
```

begin

unbundle intro-cong-syntax

lemma (in *semilattice-set*) Union: assumes finite $I I \neq \{\}$ assumes $\bigwedge i. i \in I \Longrightarrow finite (Z i)$ assumes $\bigwedge i. i \in I \Longrightarrow Z i \neq \{\}$ shows $F (\bigcup (Z `I)) = F ((\lambda i. (F (Z i))) `I)$ $\langle proof \rangle$

This is similar to the existing *hom-Max-commute* with the crucial difference that it works even if the function is a homomorphism between distinct lattices. An example application is Max (int 'A) = int (Max A).

```
lemma hom-Max-commute':
  assumes finite A \ A \neq \{\}
  assumes \bigwedge x \ y. \ x \in A \Longrightarrow y \in A \Longrightarrow max \ (f \ x) \ (f \ y) = f \ (max \ x \ y)
  shows Max (f \cdot A) = f (Max A)
  \langle proof \rangle
context inner-algorithm-fix-A
begin
definition t_c
  where t_c \ \psi \ \sigma = (Max \ ((\lambda j. \ \tau_1 \ \psi \ A \ \sigma \ j + \sigma) \ ` \{.. < b\})) - b - exp + 9
definition s_c
  where s_c \ \psi \ \sigma = nat \ (t_c \ \psi \ \sigma)
definition p_c
  where p_c \ \psi \ \sigma = card \ \{j \in \{.. < b\}. \ \tau_1 \ \psi \ A \ \sigma \ j + \sigma \ge s_c \ \psi \ \sigma\}
definition Y_c
  where Y_c \psi \sigma = 2 \hat{s}_c \psi \sigma * \varrho \text{-inv} (p_c \psi \sigma)
```

lemma s_c -eq-s: assumes $(f,g,h) \in sample-pro \Psi$ assumes $\sigma \leq s f$ shows s_c (f,g,h) $\sigma = s f$ $\langle proof \rangle$ lemma p_c -eq-p: assumes $(f,g,h) \in sample-pro \Psi$ assumes $\sigma \leq s f$ shows p_c (f,g,h) $\sigma = p$ (f,g,h) $\langle proof \rangle$ lemma Y_c -eq-Y: assumes $(f,g,h) \in sample-pro \Psi$ assumes $\sigma \leq s f$ shows Y_c (f,g,h) $\sigma = Y$ (f,g,h) $\langle proof \rangle$ **lemma** accuracy-single: measure Ψ { ψ . $\exists \sigma \leq q$ -max. | $Y_c \ \psi \ \sigma$ - real X| > $\varepsilon * X$ } $\leq 1/2^2$ (**is** $?L \leq ?R)$ $\langle proof \rangle$ **lemma** *estimate1-eq*: assumes j < lshows estimate1 ($\tau_2 \ \omega \ A \ \sigma, \sigma$) $j = Y_c \ (\omega \ j) \ \sigma \ (is \ ?L = ?R)$ $\langle proof \rangle$ **lemma** *estimate-result-1*: measure $\Omega \left\{ \omega. \left(\exists \sigma \leq q \text{-max. } \varepsilon \ast X < | \text{estimate } (\tau_2 \ \omega \ A \ \sigma, \sigma) - X | \right) \right\} \leq \delta/2 \text{ (is } ?L \leq ?R)$ $\langle proof \rangle$ theorem estimate-result: measure $\Omega \{ \omega. | estimate (\tau \ \omega \ A) - X | > \varepsilon * X \} \leq \delta$ $(\mathbf{is} ?L \leq ?R)$ $\langle proof \rangle$ end **lemma** (in inner-algorithm) estimate-result: assumes $A \subseteq \{.. < n\} A \neq \{\}$ shows measure Ω { ω . |estimate ($\tau \omega A$) - real (card A)| > $\varepsilon *$ real (card A)} $\leq \delta$ (is ?L \leq ?R) $\langle proof \rangle$

unbundle *no-intro-cong-syntax*

end

10 Outer Algorithm

This section introduces the final solution with optimal size space usage. Internally it relies on the inner algorithm described in Section 6, dependending on the parameters n, ε and δ it either uses the inner algorithm directly or if ε^{-1} is larger than $\ln n$ it runs $\frac{\varepsilon^{-1}}{\ln \ln n}$ copies of the inner algorithm (with the modified failure probability $\frac{1}{\ln n}$) using an expander to select its seeds. The theorems below verify that the probability that the relative accuracy of the median of the copies is too large is below ε . theory Distributed-Distinct-Elements-Outer-Algorithm imports Distributed-Distinct-Elements-Accuracy Prefix-Free-Code-Combinators.Prefix-Free-Code-Combinators Frequency-Moments.Landau-Ext Landau-Symbols.Landau-More

 \mathbf{begin}

unbundle intro-cong-syntax

The following are non-asymptotic hard bounds on the space usage for the sketches and seeds repsectively. The end of this section contains a proof that the sum is asymptotically in $\mathcal{O}(\ln(\varepsilon^{-1})\delta^{-1} + \ln n)$.

definition state-space-usage = $(\lambda(n,\varepsilon,\delta)$. $2^{40} * (\ln(1/\delta)+1)/\varepsilon^2 + \log 2 (\log 2 n + 3))$ definition seed-space-usage = $(\lambda(n,\varepsilon,\delta)$. $2^{30}+2^{23}*\ln n+48*(\log 2(1/\varepsilon)+16)^2+336*\ln(1/\delta))$

```
locale outer-algorithm =

fixes n :: nat

fixes \delta :: real

fixes \varepsilon :: real

assumes n-gt-0: n > 0

assumes \delta-gt-0: \delta > 0 and \delta-lt-1: \delta < 1

assumes \varepsilon-gt-0: \varepsilon > 0 and \varepsilon-lt-1: \varepsilon < 1

begin
```

```
definition n_0 where n_0 = max (real n) (exp (exp 5))
definition stage-two where stage-two = (\delta < (1/\ln n_0))
definition \delta_i :: real where \delta_i = (if \text{ stage-two then } (1/\ln n_0) \text{ else } \delta)
definition m :: nat where m = (if \text{ stage-two then nat } \lceil 4 * \ln (1/\delta)/\ln (\ln n_0) \rceil \text{ else } 1)
definition \alpha where \alpha = (if \text{ stage-two then } (1/\ln n_0) \text{ else } 1)
```

lemma *m*-lbound: assumes stage-two shows $m \ge 4 * \ln (1 / \delta) / \ln(\ln n_0)$ $\langle proof \rangle$

lemma *n*-lbound: $n_0 \ge exp \ (exp \ 5) \ ln \ n_0 \ge exp \ 5 \ 5 \le ln \ (ln \ n_0) \ ln \ n_0 > 1 \ n_0 > 1$ $\langle proof \rangle$

 $\begin{array}{l} \textbf{lemma } \delta 1 \text{-} gt \text{-} 0 \colon 0 < \delta_i \\ \langle proof \rangle \end{array}$

 $\begin{array}{l} \textbf{lemma } \delta 1 \text{-} lt \text{-} 1 \text{:} \ \delta_i < 1 \\ \langle proof \rangle \end{array}$

lemma *m*-gt-0-aux: assumes stage-two shows $1 \leq \ln (1 / \delta) / \ln (\ln n_0)$ $\langle proof \rangle$

lemma *m*-*g*t- θ : *m* > θ $\langle proof \rangle$

 $\begin{array}{ll} \textbf{lemma} \ \alpha \textit{-gt-0} \colon \alpha > \ 0 \\ \langle \textit{proof} \, \rangle \end{array}$

lemma α -le-1: $\alpha \leq 1$

 $\langle proof \rangle$

sublocale I: inner-algorithm n $\delta_i \in \langle proof \rangle$

abbreviation Θ where $\Theta \equiv \mathcal{E} \ m \ \alpha \ I.\Omega$

lemma Θ : $m > \theta \alpha > \theta \langle proof \rangle$

type-synonym state = inner-algorithm.state list

fun single :: nat \Rightarrow nat \Rightarrow state **where** single $\vartheta x = map (\lambda j. I. single (pro-select <math>\Theta \ \vartheta \ j) x) [\theta...< m]$

fun merge :: state \Rightarrow state \Rightarrow state **where** merge $x \ y = map \ (\lambda(x,y). \ I.merge \ x \ y) \ (zip \ x \ y)$

fun estimate :: state \Rightarrow real **where** estimate x = median m (λi . I.estimate (x ! i))

definition $\nu :: nat \Rightarrow nat set \Rightarrow state$ where $\nu \ \vartheta \ A = map \ (\lambda i. \ I.\tau \ (pro-select \ \Theta \ \vartheta \ i) \ A) \ [0..<m]$

The following three theorems verify the correctness of the algorithm. The term τ is a mathematical description of the sketch for a given subset, while *local.single*, *local.merge* are the actual functions that compute the sketches.

theorem merge-result: merge $(\nu \ \omega \ A) \ (\nu \ \omega \ B) = \nu \ \omega \ (A \cup B)$ (is ?L = ?R) $\langle proof \rangle$

theorem single-result: single $\omega x = \nu \omega \{x\}$ (is ?L = ?R) $\langle proof \rangle$

theorem estimate-result: **assumes** $A \subseteq \{..<n\} A \neq \{\}$ **defines** $p \equiv (pmf\text{-}of\text{-}set \{..<pro\text{-}size \Theta\})$ **shows** measure $p \{\omega. | estimate (\nu \ \omega \ A) - real (card \ A) | > \varepsilon * real (card \ A) \} \le \delta$ (is $?L \le ?R)$ $\langle proof \rangle$

The function *encode-state* can represent states as bit strings. This enables verification of the space usage.

definition encode-state where encode-state = Lf_e I.encode-state m

lemma encode-state: is-encoding encode-state $\langle proof \rangle$

lemma state-bit-count: bit-count (encode-state ($\nu \ \omega \ A$)) \leq state-space-usage (real n, ε, δ) (is ?L \leq ?R) $\langle proof \rangle$

Encoding function for the seeds which are just natural numbers smaller than pro-size Θ .

definition encode-seed where encode-seed = Nb_e (pro-size Θ)

lemma encode-seed: is-encoding encode-seed $\langle proof \rangle$

```
lemma random-bit-count:

assumes \omega < pro-size \Theta

shows bit-count (encode-seed \omega) \leq seed-space-usage (real n, \varepsilon, \delta)

(is ?L \leq ?R)

\langle proof \rangle
```

The following is an alternative form expressing the correctness and space usage theorems. If x is expression formed by *local.single* and *local.merge* operations. Then x requires *state-space-usage* (*real* n, ε, δ) bits to encode and *estimate* x approximates the count of the distinct universe elements in the expression.

For example:

estimate (local.merge (local.single ω 1) (local.merge (local.single ω 5) (local.single ω 1))) approximates the cardinality of $\{1, 5, 1\}$ i.e. 2.

datatype sketch-tree = Single nat | Merge sketch-tree sketch-tree

```
fun eval :: nat \Rightarrow sketch-tree \Rightarrow state

where

eval \ \omega \ (Single \ x) = single \ \omega \ x \mid

eval \ \omega \ (Merge \ x \ y) = merge \ (eval \ \omega \ x) \ (eval \ \omega \ y)

fun sketch-tree-set :: sketch-tree \Rightarrow nat set

where

sketch-tree-set \ (Single \ x) = \{x\} \mid

sketch-tree-set \ (Merge \ x \ y) = sketch-tree-set \ x \cup sketch-tree-set \ y

theorem correctness:

fixes X

assumes sketch-tree-set \ t \subseteq \{..<n\}

defines p \equiv pmf-of-set \{..<pro-size \ \Theta\}

defines X \equiv real \ (card \ (sketch-tree-set \ t))

shows measure p \ \{\omega. \ | estimate \ (eval \ \omega \ t) - X| > \varepsilon * X\} \le \delta \ (is \ ?L \le ?R)

\langle proof \rangle
```

theorem space-usage: **assumes** $\omega < pro-size \Theta$ **shows** *bit-count* (encode-state (eval ω t)) \leq state-space-usage (real n, ε, δ) (is ?A) *bit-count* (encode-seed ω) \leq seed-space-usage (real n, ε, δ) (is ?B) $\langle proof \rangle$

end

The functions *state-space-usage* and *seed-space-usage* are exact bounds on the space usage for the state and the seed. The following establishes asymptotic bounds with respect to the limit $n, \delta^{-1}, \varepsilon^{-1} \to \infty$.

context begin

Some local notation to ease proofs about the asymptotic space usage of the algorithm:

private definition *n*-of :: real × real × real ⇒ real where *n*-of = $(\lambda(n, \varepsilon, \delta), n)$ private definition δ -of :: real × real × real ⇒ real where δ -of = $(\lambda(n, \varepsilon, \delta), \delta)$ private definition ε -of :: real × real × real ⇒ real where ε -of = $(\lambda(n, \varepsilon, \delta), \varepsilon)$

private abbreviation $F :: (real \times real \times real)$ filter

where $F \equiv (at\text{-}top \times_F at\text{-}right \ 0 \times_F at\text{-}right \ 0)$

private lemma var-simps:

 $\begin{array}{l} n\text{-}of = fst\\ \varepsilon\text{-}of = (\lambda x. \ fst \ (snd \ x))\\ \delta\text{-}of = (\lambda x. \ snd \ (snd \ x))\\ \langle proof \rangle \ \textbf{lemma} \ evt\text{-}n: \ eventually \ (\lambda x. \ n\text{-}of \ x \ge n) \ F\\ \langle proof \rangle \ \textbf{lemma} \ evt\text{-}n\text{-}1: \ \forall \ F \ x \ in \ F. \ 0 \le ln \ (n\text{-}of \ x))\\ \langle proof \rangle \ \textbf{lemma} \ evt\text{-}n\text{-}2: \ \forall \ F \ x \ in \ F. \ 0 \le ln \ (ln \ (n\text{-}of \ x))\\ \langle proof \rangle \ \textbf{lemma} \ evt\text{-}\varepsilon\text{-} i: \ eventually \ (\lambda x. \ 1/\varepsilon\text{-}of \ x \ge \varepsilon \land \varepsilon\text{-}of \ x > 0) \ F\\ \langle proof \rangle \ \textbf{lemma} \ evt\text{-}\delta\text{-}1: \ \forall \ F \ x \ in \ F. \ 0 \le ln \ (1 \ / \ \delta\text{-}of \ x)\\ \langle proof \rangle \ \textbf{lemma} \ evt\text{-}\delta\text{-}1: \ \forall \ F \ x \ in \ F. \ 0 \le ln \ (1 \ / \ \delta\text{-}of \ x)\\ \langle proof \rangle \end{array}$

theorem asymptotic-state-space-complexity: state-space-usage $\in O[F](\lambda(n, \varepsilon, \delta). \ln (1/\delta)/\varepsilon^2 + \ln (\ln n))$ (is $- \in O[?F](?rhs))$ $\langle proof \rangle$

theorem asymptotic-seed-space-complexity:

seed-space-usage $\in O[F](\lambda(n, \varepsilon, \delta))$. $\ln(1/\delta) + \ln(1/\varepsilon)^2 + \ln n$ (is $- \in O[?F](?rhs))$ (proof)

definition space-usage x = state-space-usage x + seed-space-usage x

```
theorem asymptotic-space-complexity:
space-usage \in O[at\text{-top } \times_F at\text{-right } 0 \times_F at\text{-right } 0](\lambda(n, \varepsilon, \delta). \ln(1/\delta)/\varepsilon^2 + \ln n) \langle proof \rangle
```

 \mathbf{end}

unbundle no-intro-cong-syntax

end

References

- 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 *Randomization and Approximation Techniques in Computer Science*, pages 1–10. Springer Berlin Heidelberg, 2002.
- [3] J. Błasiok. Optimal streaming and tracking distinct elements with high probability. ACM Trans. Algorithms, 16(1):3:1–3:28, 2020.
- [4] P. Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. Journal of Computer and System Sciences, 31(2):182–209, 1985.
- [5] P. B. Gibbons and S. Tirthapura. Estimating simple functions on the union of data streams. In Proceedings of the Thirteenth Annual ACM Symposium on Parallel Algorithms and Architectures, SPAA '01, pages 281–291, 2001.
- [6] V. Guruswami, C. Umans, and S. Vadhan. Unbalanced expanders and randomness extractors from parvaresh-vardy codes. J. ACM, 56(4), jul 2009.
- [7] D. M. Kane, J. Nelson, and D. P. Woodruff. An optimal algorithm for the distinct elements problem. In Proceedings of the Twenty-Ninth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS '10, pages 41–52, New York, 2010.

- [8] E. Karayel. Finite fields. Archive of Formal Proofs, June 2022. https://isa-afp.org/entries/ Finite_Fields.html, Formal proof development.
- [9] E. Karayel. Formalization of randomized approximation algorithms for frequency moments. *Archive of Formal Proofs*, April 2022. https://isa-afp.org/entries/Frequency_Moments.html, Formal proof development.
- [10] E. Karayel. Expander graphs. Archive of Formal Proofs, March 2023. https://isa-afp.org/ entries/Expander_Graphs.html, Formal proof development.
- [11] D. Woodruff. Optimal space lower bounds for all frequency moments. In Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '04, pages 167–175, USA, 2004. Society for Industrial and Applied Mathematics.