

The Median Method

Emin Karayel

July 1, 2022

Abstract

The median method is an amplification result for randomized approximation algorithms described in [1]. Given an algorithm whose result is in a desired interval with a probability larger than $\frac{1}{2}$, it is possible to improve the success probability, by running the algorithm multiple times independently and using the median. In contrast to using the mean, the amplification of the success probability grows exponentially with the number of independent runs.

This entry contains a formalization of the underlying theorem: Given a sequence of n independent random variables, which are in a desired interval with a probability $\frac{1}{2} + \alpha$. Then their median will be in the desired interval with a probability of $1 - \exp(-2\alpha^2 n)$. In particular, the success probability approaches 1 exponentially with the number of variables.

In addition to that, this entry also contains a proof that order-statistics of Borel-measurable random variables are themselves measurable and that generalized intervals in linearly ordered Borel-spaces are measurable.

Contents

1	Intervals are Borel measurable	1
2	Order statistics are Borel measurable	4
3	The Median Method	8
4	Some additional results about the median	14

1 Intervals are Borel measurable

theory *Median*

imports *HOL-Probability.Hoeffding HOL-Library.Multiset*

begin

This section contains a proof that intervals are Borel measurable, where an interval is defined as a convex subset of linearly ordered space, more precisely, a set is an interval, if for each triple of points $x < y < z$: If x and z are in the set so is y . This includes ordinary intervals like $\{a..b\}$, $\{a<..**b\}**$ but also for example $\{x::rat. x * x < (2::rat)\}$ which cannot be expressed in the standard notation.

In the *HOL-Analysis.Borel-Space* there are proofs for the measurability of each specific type of interval, but those unfortunately do not help if we want to express the result about the median bound for arbitrary types of intervals.

definition *interval* :: ('a :: linorder) set \Rightarrow bool **where**
interval I = ($\forall x y z. x \in I \longrightarrow z \in I \longrightarrow x \leq y \longrightarrow y \leq z \longrightarrow y \in I$)

definition *up-ray* :: ('a :: linorder) set \Rightarrow bool **where**
up-ray I = ($\forall x y. x \in I \longrightarrow x \leq y \longrightarrow y \in I$)

lemma *up-ray-borel*:
assumes *up-ray* (I :: (('a :: linorder-topology) set))
shows I \in borel

proof (cases closed I)
case True
then show ?thesis **using** borel-closed **by** blast
next
case False
hence b: \neg closed I **by** blast

have open I
proof (rule Topological-Spaces.openI)
fix x
assume c:x \in I
show $\exists T. open T \wedge x \in T \wedge T \subseteq I$
proof (cases $\exists y. y < x \wedge y \in I$)
case True
then obtain y **where** a:y < x \wedge y \in I **by** blast
have open {y<..} **by** simp
moreover have x \in {y<..} **using** a **by** simp
moreover have {y<..} \subseteq I
apply (rule subsetI)
using a *assms*(1) **apply** (simp add: up-ray-def)
by (metis less-le-not-le)
ultimately show ?thesis **by** blast
next
case False
hence I \subseteq {x..} **using** linorder-not-less **by** auto
moreover have {x..} \subseteq I
using c *assms*(1) **apply** (simp add: up-ray-def)
by blast
ultimately have I = {x..}
by (rule order-antisym)

moreover have $\text{closed } \{x..\}$ **by** *simp*
ultimately have False **using** b **by** *auto*
then show $?thesis$ **by** *simp*
qed
qed
then show $?thesis$ **by** *simp*
qed

definition $\text{down-ray} :: ('a :: \text{linorder}) \text{ set} \Rightarrow \text{bool}$ **where**
 $\text{down-ray } I = (\forall x y. y \in I \longrightarrow x \leq y \longrightarrow x \in I)$

lemma down-ray-borel :
assumes $\text{down-ray } (I :: (('a :: \text{linorder-topology}) \text{ set}))$
shows $I \in \text{borel}$
proof –
have $\text{up-ray } (-I)$ **using** *assms*
by (*simp add: up-ray-def down-ray-def, blast*)
hence $(-I) \in \text{borel}$ **using** up-ray-borel **by** *blast*
thus $I \in \text{borel}$
by (*metis borel-comp double-complement*)
qed

Main result of this section:

lemma interval-borel :
assumes $\text{interval } (I :: (('a :: \text{linorder-topology}) \text{ set}))$
shows $I \in \text{borel}$
proof (*cases* $I = \{\}$)
case True
then show $?thesis$ **by** *simp*
next
case False
then obtain x **where** $a : x \in I$ **by** *blast*
have $\bigwedge y z. y \in I \cup \{x..\} \Longrightarrow y \leq z \Longrightarrow z \in I \cup \{x..\}$
by (*metis assms a interval-def IntE UnE Un-Int-eq(1) Un-Int-eq(2) atLeast-iff nle-le order.trans*)
hence $\text{up-ray } (I \cup \{x..\})$
using up-ray-def **by** *blast*
hence $b : I \cup \{x..\} \in \text{borel}$
using up-ray-borel **by** *blast*

have $\bigwedge y z. y \in I \cup \{..x\} \Longrightarrow z \leq y \Longrightarrow z \in I \cup \{..x\}$
by (*metis assms a interval-def UnE UnI1 UnI2 atMost-iff dual-order.trans linorder-le-cases*)
hence $\text{down-ray } (I \cup \{..x\})$
using down-ray-def **by** *blast*
hence $c : I \cup \{..x\} \in \text{borel}$
using down-ray-borel **by** *blast*

have $I = (I \cup \{x..\}) \cap (I \cup \{..x\})$

```

    using a by fastforce

  then show ?thesis using b c
    by (metis sets.Int)
qed

```

2 Order statistics are Borel measurable

This section contains a proof that order statistics of Borel measurable random variables are themselves Borel measurable.

The proof relies on the existence of branch-free comparison-sort algorithms. Given a sequence length these algorithms perform compare-swap operations on predefined pairs of positions. In particular the result of a comparison does not affect future operations. An example for a branch-free comparison sort algorithm is shell-sort and also bubble-sort without the early exit.

The advantage of using such a comparison-sort algorithm is that it can be lifted to work on random variables, where the result of a comparison-swap operation on two random variables X and Y can be represented as the expressions $\lambda\omega. \min (X \ \omega) (Y \ \omega)$ and $\lambda\omega. \max (X \ \omega) (Y \ \omega)$.

Because taking the point-wise minimum (resp. maximum) of two random variables is still Borel measurable, and because the entire sorting operation can be represented using such compare-swap operations, we can show that all order statistics are Borel measurable.

fun *sort-primitive* **where**

```

  sort-primitive i j f k = (if k = i then min (f i) (f j) else (if k = j then max (f i)
(f j) else f k))

```

fun *sort-map* **where**

```

  sort-map f n = fold id [sort-primitive j i. i <- [0..<n], j <- [0..<i]] f

```

lemma *sort-map-ind*:

```

  sort-map f (Suc n) = fold id [sort-primitive j n. j <- [0..<n]] (sort-map f n)
  by simp

```

lemma *sort-map-strict-mono*:

```

  fixes f :: nat => 'b :: linorder

```

```

  shows j < n => i < j => sort-map f n i <= sort-map f n j

```

proof (*induction n arbitrary: i j*)

```

  case 0

```

```

  then show ?case by simp

```

next

```

  case (Suc n)

```

```

  define g where g = ( $\lambda k. \text{fold id [sort-primitive j n. j <- [0..<k]] (sort-map f n)}$ )

```

```

  define k where k = n

```

```

  have a: ( $\forall i j. j < n \longrightarrow i < j \longrightarrow g k i \leq g k j$ )  $\wedge$  ( $\forall l. l < k \longrightarrow g k l \leq g k n$ )

```

```

proof (induction k)
  case 0
  then show ?case using Suc by (simp add:g-def del:sort-map.simps)
next
  case (Suc k)
  have g (Suc k) = sort-primitive k n (g k)
  by (simp add:g-def)
  then show ?case using Suc
  apply (cases g k k ≤ g k n)
  apply (simp add:min-def max-def)
  using less-antisym apply blast
  apply (cases g k n ≤ g k k)
  apply (simp add:min-def max-def)
  apply (metis less-antisym max.coboundedI2 max.orderE)
  by simp
qed

hence  $\bigwedge i j. j < \text{Suc } n \implies i < j \implies g \ n \ i \leq g \ n \ j$ 
  apply (simp add:k-def) using less-antisym by blast
moreover have sort-map f (Suc n) = g n
  by (simp add:sort-map-ind g-def del:sort-map.simps)
ultimately show ?case
  apply (simp del:sort-map.simps)
  using Suc by blast
qed

lemma sort-map-mono:
  fixes f :: nat  $\Rightarrow$  'b :: linorder
  shows j < n  $\implies$  i ≤ j  $\implies$  sort-map f n i ≤ sort-map f n j
  by (metis sort-map-strict-mono eq-iff le-imp-less-or-eq)

lemma sort-map-perm:
  fixes f :: nat  $\Rightarrow$  'b :: linorder
  shows image-mset (sort-map f n) (mset [0..proof -
  define is-swap where is-swap = (λ(ts :: ((nat  $\Rightarrow$  'b)  $\Rightarrow$  nat  $\Rightarrow$  'b)).  $\exists i < n. \exists j < n. ts = \text{sort-primitive } i \ j$ )
  define t :: ((nat  $\Rightarrow$  'b)  $\Rightarrow$  nat  $\Rightarrow$  'b) list
  where t = [sort-primitive j i. i <- [0..have a:  $\bigwedge x f. \text{is-swap } x \implies \text{image-mset } (x \ f) \ (\text{mset-set } \{0..
  proof -
  fix x
  fix f :: nat  $\Rightarrow$  'b :: linorder
  assume is-swap x
  then obtain i j where x-def: x = sort-primitive i j and i-bound: i < n and j-bound:j < n
  using is-swap-def by blast$ 
```

```

define inv where inv = mset-set {k. k < n ∧ k ≠ i ∧ k ≠ j}
have b:{0..<n} = {k. k < n ∧ k ≠ i ∧ k ≠ j} ∪ {i,j}
  apply (rule order-antisym, rule subsetI, simp, blast, rule subsetI, simp)
  using i-bound j-bound by meson
have c:∧k. k ∈# inv ⇒ (x f) k = f k
  by (simp add:x-def inv-def)
have image-mset (x f) inv = image-mset f inv
  apply (rule multiset-eqI)
  using c multiset.map-cong0 by force
moreover have image-mset (x f) (mset-set {i,j}) = image-mset f (mset-set
{i,j})
  apply (cases i = j)
  by (simp add:x-def max-def min-def)+
moreover have mset-set {0..<n} = inv + mset-set {i,j}
  by (simp only:inv-def b, rule mset-set-Union, simp, simp, simp)
ultimately show image-mset (x f) (mset-set {0..<n}) = image-mset f (mset-set
{0..<n})
  by simp
qed

```

```

have (∀x ∈ set t. is-swap x) ⇒ image-mset (fold id t f) (mset [0..<n]) =
image-mset f (mset [0..<n])
  by (induction t arbitrary:f, simp, simp add:a)
moreover have ∧x. x ∈ set t ⇒ is-swap x
  apply (simp add:t-def is-swap-def)
  by (meson atLeastLessThan-iff imageE less-imp-le less-le-trans)
ultimately have image-mset (fold id t f) (mset [0..<n]) = image-mset f (mset
[0..<n]) by blast
  then show ?thesis by (simp add:t-def)
qed

```

```

lemma list-eq-iff:
assumes mset xs = mset ys
assumes sorted xs
assumes sorted ys
shows xs = ys
using assms properties-for-sort by blast

```

```

lemma sort-map-eq-sort:
fixes f :: nat ⇒ ('b :: linorder)
shows map (sort-map f n) [0..<n] = sort (map f [0..<n]) (is ?A = ?B)
proof –
have mset ?A = mset ?B
  using sort-map-perm[where f=f and n=n]
  by (simp del:sort-map.simps)
moreover have sorted ?B
  by simp
moreover have sorted ?A
  apply (subst sorted-wrt-iff-nth-less)

```

```

    apply (simp del:sort-map.simps)
    by (metis sort-map-mono nat-less-le)
  ultimately show ?A = ?B
    using list-eq-iff by blast
qed

lemma order-statistics-measurable-aux:
  fixes X :: nat ⇒ 'a ⇒ ('b :: {linorder-topology, second-countable-topology})
  assumes n ≥ 1
  assumes j < n
  assumes ⋀i. i < n ⇒ X i ∈ measurable M borel
  shows (λx. (sort-map (λi. X i x) n) j) ∈ measurable M borel
proof -
  have n-ge-0:n > 0 using assms by simp
  define is-swap where is-swap = (λ(ts :: ((nat ⇒ 'b) ⇒ nat ⇒ 'b)). ∃ i < n. ∃ j
< n. ts = sort-primitive i j)
  define t :: ((nat ⇒ 'b) ⇒ nat ⇒ 'b) list
    where t = [sort-primitive j i. i <- [0..<n], j <- [0..<i]]

  define meas-ptw :: (nat ⇒ 'a ⇒ 'b) ⇒ bool
    where meas-ptw = (λf. (∀ k. k < n → f k ∈ borel-measurable M))

  have ind-step:
    ⋀x (g :: nat ⇒ 'a ⇒ 'b). meas-ptw g ⇒ is-swap x ⇒ meas-ptw (λk ω. x (λi.
g i ω) k)
  proof -
    fix x g
    assume meas-ptw g
    hence a:⋀k. k < n ⇒ g k ∈ borel-measurable M by (simp add:meas-ptw-def)
    assume is-swap x
    then obtain i j where x-def:x=sort-primitive i j and i-le:i < n and j-le:j <
n
      by (simp add:is-swap-def, blast)
    have ⋀k. k < n ⇒ (λω. x (λi. g i ω) k) ∈ borel-measurable M
  proof -
    fix k
    assume k < n
    thus (λω. x (λi. g i ω) k) ∈ borel-measurable M
      apply (simp add:x-def)
      apply (cases k = i, simp)
      using a i-le j-le borel-measurable-min apply blast
      apply (cases k = j, simp)
      using a i-le j-le borel-measurable-max apply blast
      using a by simp
    qed
  thus meas-ptw (λk ω. x (λi. g i ω) k)
    by (simp add:meas-ptw-def)
qed

```

```

have ( $\forall x \in \text{set } t. \text{is-swap } x \implies \text{meas-ptw } (\lambda k \omega. (\text{fold id } t (\lambda k. X k \omega)) k)$ )
proof (induction t rule:rev-induct)
  case Nil
  then show ?case using assms by (simp add:meas-ptw-def)
next
  case (snoc x xs)
  have a:meas-ptw ( $\lambda k \omega. \text{fold } (\lambda a. a) \text{ xs } (\lambda k. X k \omega) k$ ) using snoc by simp
  have b:is-swap x using snoc by simp
  show ?case using ind-step[OF a b] by simp
qed
moreover have  $\bigwedge x. x \in \text{set } t \implies \text{is-swap } x$ 
  apply (simp add:t-def is-swap-def)
  by (meson atLeastLessThan-iff imageE less-imp-le less-le-trans)
ultimately show ?thesis using assms
  by (simp add:t-def[symmetric] meas-ptw-def)
qed

```

Main results of this section:

lemma *order-statistics-measurable*:

```

fixes X :: nat  $\Rightarrow$  'a  $\Rightarrow$  ('b :: {linorder-topology, second-countable-topology})
assumes  $n \geq 1$ 
assumes  $j < n$ 
assumes  $\bigwedge i. i < n \implies X i \in \text{measurable } M \text{ borel}$ 
shows  $(\lambda x. (\text{sort } (\text{map } (\lambda i. X i x) [0..<n])) ! j) \in \text{measurable } M \text{ borel}$ 
apply (subst sort-map-eq-sort[symmetric])
using assms by (simp add:order-statistics-measurable-aux del:sort-map.simps)

```

definition *median where*

```

median n f = sort (map f [0..<n]) ! (n div 2)

```

lemma *median-measurable*:

```

fixes X :: nat  $\Rightarrow$  'a  $\Rightarrow$  ('b :: {linorder-topology, second-countable-topology})
assumes  $n \geq 1$ 
assumes  $\bigwedge i. i < n \implies X i \in \text{measurable } M \text{ borel}$ 
shows  $(\lambda x. \text{median } n (\lambda i. X i x)) \in \text{measurable } M \text{ borel}$ 
apply (simp add:median-def)
apply (rule order-statistics-measurable[OF assms(1) - assms(2)])
using assms(1) by force+

```

3 The Median Method

This section contains the proof for the probability that the median of independent random variables will be in an interval with high probability if the individual variables are in the same interval with probability larger than $\frac{1}{2}$. The proof starts with the elementary observation that the median of a sequence with n elements is in an interval I if at least half of them are in I . This works because after sorting the sequence the elements that will be in

the interval must necessarily form a consecutive subsequence, if its length is larger than $\frac{n}{2}$ the median must be in it.

The remainder follows the proof in [1, §2.1] using the Hoeffding inequality to estimate the probability that at least half of the sequence elements will be in the interval I .

```

lemma interval-rule:
  assumes interval I
  assumes  $a \leq x \ x \leq b$ 
  assumes  $a \in I$ 
  assumes  $b \in I$ 
  shows  $x \in I$ 
  using assms(1) apply (simp add: interval-def)
  using assms by blast

```

```

lemma sorted-int:
  assumes interval I
  assumes sorted xs
  assumes  $k < \text{length } xs \ i \leq j \ j \leq k$ 
  assumes  $xs ! i \in I \ xs ! k \in I$ 
  shows  $xs ! j \in I$ 
  apply (rule interval-rule[where  $a=xs ! i$  and  $b=xs ! k$ ])
  using assms by (simp add: sorted-nth-mono)+

```

```

lemma mid-in-interval:
  assumes  $2 * \text{length} (\text{filter } (\lambda x. x \in I) \ xs) > \text{length } xs$ 
  assumes interval I
  assumes sorted xs
  shows  $xs ! (\text{length } xs \ \text{div } 2) \in I$ 
proof –
  have  $\text{length} (\text{filter } (\lambda x. x \in I) \ xs) > 0$  using assms(1) by linarith
  then obtain  $v$  where  $v-1: v < \text{length } xs$  and  $v-2: xs ! v \in I$ 
  by (metis filter-False in-set-conv-nth length-greater-0-conv)

```

```

define  $J$  where  $J = \{k. k < \text{length } xs \wedge xs ! k \in I\}$ 

```

```

have card-J-min:  $2 * \text{card } J > \text{length } xs$ 
  using assms(1) by (simp add: J-def length-filter-conv-card)

```

```

consider
  (a)  $xs ! (\text{length } xs \ \text{div } 2) \in I \mid$ 
  (b)  $xs ! (\text{length } xs \ \text{div } 2) \notin I \wedge v > (\text{length } xs \ \text{div } 2) \mid$ 
  (c)  $xs ! (\text{length } xs \ \text{div } 2) \notin I \wedge v < (\text{length } xs \ \text{div } 2) \mid$ 
  by (metis linorder-cases v-2)
thus ?thesis
proof (cases)
  case  $a$ 
  then show ?thesis by simp
next

```

```

case b
have p:  $\bigwedge k. k \leq \text{length } xs \text{ div } 2 \implies xs ! k \notin I$ 
  using b v-2 sorted-int[OF assms(2) assms(3) v-1, where j=length xs div 2]
apply simp by blast
have card J  $\leq$  card {Suc (length xs div 2)..<length xs}
  apply (rule card-mono, simp)
  apply (rule subsetI, simp add:J-def not-less-eq-eq[symmetric])
  using p by metis
hence card J  $\leq$  length xs - (Suc (length xs div 2))
  using card-atLeastLessThan by metis
hence length xs  $\leq$  2*(length xs - (Suc (length xs div 2)))
  using card-J-min by linarith
hence False
  apply (simp add:nat-distrib)
  apply (subst (asm) le-diff-conv2) using b v-1 apply linarith
  by simp
then show ?thesis by simp
next
case c
have p:  $\bigwedge k. k \geq \text{length } xs \text{ div } 2 \implies k < \text{length } xs \implies xs ! k \notin I$ 
  using c v-1 v-2 sorted-int[OF assms(2) assms(3), where i=v and j=length
xs div 2] apply simp by blast
have card J  $\leq$  card {0..<(length xs div 2)}
  apply (rule card-mono, simp)
  apply (rule subsetI, simp add:J-def not-less-eq-eq[symmetric])
  using p linorder-le-less-linear by blast
hence card J  $\leq$  (length xs div 2)
  using card-atLeastLessThan by simp
then show ?thesis using card-J-min by linarith
qed
qed

```

lemma median-est:

```

assumes interval I
assumes 2*card {k. k < n  $\wedge$  f k  $\in$  I} > n
shows median n f  $\in$  I
proof -
have a: {k. k < n  $\wedge$  f k  $\in$  I} = {i. i < n  $\wedge$  map f [0..<n] ! i  $\in$  I}
  apply (rule order-antisym, rule subsetI, simp)
  by (rule subsetI, simp, metis add-0 diff-zero nth-map-upt)

```

```

show ?thesis
  apply (simp add:median-def)
  apply (rule mid-in-interval[where I=I and xs=sort (map f [0..<n]), simpli-
fied])
  using assms a apply (simp add:filter-sort comp-def length-filter-conv-card)
  by (simp add:assms)
qed

```

Main results of this section:

theorem (in *prob-space*) *median-bound*:

fixes $n :: \text{nat}$

fixes $I :: ('b :: \{\text{linorder-topology, second-countable-topology}\}) \text{ set}$

assumes *interval* I

assumes $\alpha > 0$

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

assumes *indep-vars* $(\lambda-. \text{borel}) X \{0..<n\}$

assumes $n \geq -\ln \varepsilon / (2 * \alpha^2)$

assumes $\bigwedge i. i < n \implies \mathcal{P}(\omega \text{ in } M. X i \omega \in I) \geq 1/2 + \alpha$

shows $\mathcal{P}(\omega \text{ in } M. \text{median } n (\lambda i. X i \omega) \in I) \geq 1 - \varepsilon$

proof –

define $Y :: \text{nat} \Rightarrow 'a \Rightarrow \text{real}$ **where** $Y = (\lambda i. \text{indicator } I \circ (X i))$

define t **where** $t = (\sum i = 0..<n. \text{expectation } (Y i)) - n/2$

have $0 < -\ln \varepsilon / (2 * \alpha^2)$

apply (*rule divide-pos-pos*)

apply (*simp, subst ln-less-zero-iff*)

using *assms* **by** *auto*

also have $\dots \leq \text{real } n$ **using** *assms* **by** *simp*

finally have $\text{real } n > 0$ **by** *simp*

hence $n - \text{ge} - 1 : n \geq 1$ **by** *linarith*

hence $n - \text{ge} - 0 : n > 0$ **by** *simp*

have *ind-comp*: $\bigwedge i. \text{indicator } I \circ (X i) = \text{indicator } \{\omega. X i \omega \in I\}$

by (*rule ext, simp add:indicator-def comp-def*)

have $\alpha * n \leq (\sum i = 0..<n. 1/2 + \alpha) - n/2$

by (*simp add:algebra-simps*)

also have $\dots \leq (\sum i = 0..<n. \text{expectation } (Y i)) - n/2$

apply (*rule diff-right-mono, rule sum-mono*)

using *assms*(6) **by** (*simp add:Y-def ind-comp Collect-conj-eq inf-commute*)

also have $\dots = t$ **by** (*simp add:t-def*)

finally have $t - \text{ge} - a : t \geq \alpha * n$ **by** *simp*

have $d : 0 \leq \alpha * n$

apply (*rule mult-nonneg-nonneg*)

using *assms*(2) $n - \text{ge} - 0$ **by** *simp+*

also have $\dots \leq t$ **using** $t - \text{ge} - a$ **by** *simp*

finally have $t - \text{ge} - 0 : t \geq 0$ **by** *simp*

have $(\alpha * n)^2 \leq t^2$ **using** $t - \text{ge} - a$ d *power-mono* **by** *blast*

hence $t - \text{ge} - a - \text{sq} : \alpha^2 * \text{real } n * \text{real } n \leq t^2$

by (*simp add:algebra-simps power2-eq-square*)

have $Y - \text{indep} : \text{indep-vars } (\lambda-. \text{borel}) Y \{0..<n\}$

apply (*subst Y-def*)

apply (*rule indep-vars-compose*[**where** $M' = (\lambda-. \text{borel}), OF \text{ assms}(4)$])

using *interval-borel*[$OF \text{ assms}(1)$] **by** *simp*

hence $b:\text{Hoeffding-ineq } M \{0..<n\} Y (\lambda i. 0) (\lambda i. 1)$
apply (*simp add:Hoeffding-ineq-def indep-interval-bounded-random-variables-def*)
by (*simp add:prob-space-axioms indep-interval-bounded-random-variables-axioms-def*
Y-def Y-indep)

have $c: \bigwedge \omega. (\sum i = 0..<n. Y i \omega) > n/2 \implies \text{median } n (\lambda i. X i \omega) \in I$
proof –

fix ω
assume $(\sum i = 0..<n. Y i \omega) > n/2$
hence $n < 2 * \text{card} (\{0..<n\} \cap \{i. X i \omega \in I\})$
by (*simp add:Y-def indicator-def*)
also have $\dots = 2 * \text{card} \{i. i < n \wedge X i \omega \in I\}$
apply (*simp, rule arg-cong[where f=card]*)
by (*rule order-antisym, rule subsetI, simp, rule subsetI, simp*)
finally have $2 * \text{card} \{i. i < n \wedge X i \omega \in I\} > n$ **by** *simp*
thus $\text{median } n (\lambda i. X i \omega) \in I$
using *median-est[OF assms(1)]* **by** *simp*

qed

have $1 - \varepsilon \leq 1 - \exp(- (2 * \alpha^2 * \text{real } n))$
apply (*simp, subst ln-ge-iff[symmetric]*)
using *assms(3)* **apply** *simp*
using *assms(5)* **apply** (*subst (asm) pos-divide-le-eq*)
apply (*simp add: assms(2) power2-eq-square*)
by (*simp add: mult-of-nat-commute*)
also have $\dots \leq 1 - \exp(- (2 * t^2 / \text{real } n))$
apply *simp*
apply (*subst pos-le-divide-eq*) **using** *n-ge-0* **apply** *simp*
using *t-ge-a-sq* **by** *linarith*
also have $\dots \leq 1 - \mathcal{P}(\omega \text{ in } M. (\sum i = 0..<n. Y i \omega) \leq n/2)$
using *Hoeffding-ineq.Hoeffding-ineq-le[OF b, where $\varepsilon=t$, simplified]* *n-ge-0*

t-ge-0

by (*simp add:t-def*)
also have $\dots = \mathcal{P}(\omega \text{ in } M. (\sum i = 0..<n. Y i \omega) > n/2)$
apply (*subst prob-compl[symmetric]*)
apply *measurable*
using *Y-indep* **apply** (*simp add:indep-vars-def*)
apply (*rule arg-cong2[where f=measure], simp*)
by (*rule order-antisym, rule subsetI, simp add:not-le, rule subsetI, simp add:not-le*)
also have $\dots \leq \mathcal{P}(\omega \text{ in } M. \text{median } n (\lambda i. X i \omega) \in I)$
apply (*rule finite-measure-mono*)
apply (*rule subsetI*) **using** *c* **apply** *simp*
using *interval-borel[OF assms(1)]* **apply** *measurable*
apply (*rule median-measurable[OF n-ge-1]*)
using *assms(4)* **by** (*simp add:indep-vars-def*)
finally show *?thesis* **by** *simp*

qed

This is a specialization of the above to closed real intervals.

corollary (in prob-space) median-bound-1:

assumes $\alpha > 0$
assumes $\varepsilon \in \{0 < \cdot < 1\}$
assumes *indep-vars* (λ -. *borel*) $X \{0..<n\}$
assumes $n \geq -\ln \varepsilon / (2 * \alpha^2)$
assumes $\forall i \in \{0..<n\}. \mathcal{P}(\omega \text{ in } M. X i \omega \in (\{a..b\} :: \text{real set})) \geq 1/2 + \alpha$
shows $\mathcal{P}(\omega \text{ in } M. \text{median } n (\lambda i. X i \omega) \in \{a..b\}) \geq 1 - \varepsilon$
apply (*rule median-bound*[*OF - assms(1) assms(2) assms(3) assms(4)*])
apply (*simp add:interval-def*)
using *assms(5)* **by** *auto*

This is a specialization of the above, where $\alpha = \frac{1}{6}$ and the interval is described using a mid point μ and radius δ . The choice of $\alpha = \frac{1}{6}$ implies a success probability per random variable of $\frac{2}{3}$. It is a commonly chosen success probability for Monte-Carlo algorithms (cf. [2, §4] or [3, §1]).

corollary (in prob-space) median-bound-2:

fixes $\mu \delta :: \text{real}$
assumes $\varepsilon \in \{0 < \cdot < 1\}$
assumes *indep-vars* (λ -. *borel*) $X \{0..<n\}$
assumes $n \geq -18 * \ln \varepsilon$
assumes $\bigwedge i. i < n \implies \mathcal{P}(\omega \text{ in } M. \text{abs } (X i \omega - \mu) > \delta) \leq 1/3$
shows $\mathcal{P}(\omega \text{ in } M. \text{abs } (\text{median } n (\lambda i. X i \omega) - \mu) \leq \delta) \geq 1 - \varepsilon$

proof –

have $b: \bigwedge i. i < n \implies \text{space } M - \{\omega \in \text{space } M. X i \omega \in \{\mu - \delta.. \mu + \delta\}\} = \{\omega \in \text{space } M. \text{abs } (X i \omega - \mu) > \delta\}$

apply (*rule order-antisym, rule subsetI, simp, linarith*)
by (*rule subsetI, simp, linarith*)

have $\bigwedge i. i < n \implies 1 - \mathcal{P}(\omega \text{ in } M. X i \omega \in \{\mu - \delta.. \mu + \delta\}) \leq 1/3$

apply (*subst prob-compl[symmetric]*)
apply (*measurable*)
using *assms(2)* **apply** (*simp add:indep-vars-def*)
apply (*subst b, simp*)
using *assms(4)* **by** *simp*

hence $a: \bigwedge i. i < n \implies \mathcal{P}(\omega \text{ in } M. X i \omega \in \{\mu - \delta.. \mu + \delta\}) \geq 2/3$ **by** *simp*

have $1 - \varepsilon \leq \mathcal{P}(\omega \text{ in } M. \text{median } n (\lambda i. X i \omega) \in \{\mu - \delta.. \mu + \delta\})$

apply (*rule median-bound-1*[*OF - assms(1) assms(2), where $\alpha=1/6$*], *simp*)
using *assms(3)* **apply** (*simp add:power2-eq-square*)
using *a* **by** *simp*

also have $\dots = \mathcal{P}(\omega \text{ in } M. \text{abs } (\text{median } n (\lambda i. X i \omega) - \mu) \leq \delta)$

apply (*rule arg-cong2*[**where** *f=measure*], *simp*)
apply (*rule order-antisym, rule subsetI, simp, linarith*)
by (*rule subsetI, simp, linarith*)

finally show *?thesis* **by** *simp*

qed

4 Some additional results about the median

lemma *sorted-mono-map*:
 assumes *sorted xs*
 assumes *mono f*
 shows *sorted (map f xs)*
 using *assms* **apply** (*simp add:sorted-wrt-map*)
 apply (*rule sorted-wrt-mono-rel[where P=(≤)]*)
 by (*simp add:mono-def, simp*)

This could be added to *HOL.List*:

lemma *map-sort*:
 assumes *mono f*
 shows *sort (map f xs) = map f (sort xs)*
 apply (*rule properties-for-sort*)
 apply *simp*
 by (*rule sorted-mono-map, simp, simp add:assms*)

lemma *median-cong*:
 assumes $\bigwedge i. i < n \implies f i = g i$
 shows *median n f = median n g*
 apply (*cases n = 0, simp add:median-def*)
 apply (*simp add:median-def*)
 apply (*rule arg-cong2[where f=(!)]*)
 apply (*rule arg-cong[where f=sort], rule map-cong, simp, simp add:assms*)
 by *simp*

lemma *median-restrict*:
 median n ($\lambda i \in \{0..<n\}.f i$) = median n f
 by (*rule median-cong, simp*)

lemma *median-commute-mono*:
 assumes $n > 0$
 assumes *mono g*
 shows $g (\text{median } n \ f) = \text{median } n \ (g \circ f)$
 apply (*simp add: median-def del:map-map*)
 apply (*subst map-map[symmetric]*)
 apply (*subst map-sort[OF assms(2)]*)
 apply (*subst nth-map, simp*) **using** *assms* **apply** *fastforce*
 by *simp*

lemma *median-rat*:
 assumes $n > 0$
 shows *real-of-rat (median n f) = median n ($\lambda i. \text{real-of-rat } (f i)$)*
 apply (*subst (2) comp-def[where g=f, symmetric]*)
 apply (*rule median-commute-mono[OF assms(1)]*)
 by (*simp add: mono-def of-rat-less-eq*)

lemma *median-const*:

```

assumes  $k > 0$ 
shows  $\text{median } k (\lambda i \in \{0..<k\}. a) = a$ 
proof –
  have  $b: \text{sorted } (\text{map } (\lambda-. a) [0..<k])$ 
    by ( $\text{subst sorted-wrt-map, simp}$ )
  have  $a: \text{sort } (\text{map } (\lambda-. a) [0..<k]) = \text{map } (\lambda-. a) [0..<k]$ 
    by ( $\text{subst sorted-sort-id[OF } b], \text{simp}$ )
  have  $\text{median } k (\lambda i \in \{0..<k\}. a) = \text{median } k (\lambda-. a)$ 
    by ( $\text{subst median-restrict, simp}$ )
  also have  $\dots = a$ 
    apply ( $\text{simp add:median-def } a$ )
    apply ( $\text{subst nth-map}$ )
    using  $\text{assms}$  by  $\text{simp+}$ 
  finally show  $?thesis$  by  $\text{simp}$ 
qed

end

```

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 *Randomization and Approximation Techniques in Computer Science*, pages 1–10. Springer Berlin Heidelberg, 2002.
- [3] 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.