

Universal Hash Families

Emin Karayel

March 24, 2023

Abstract

A k -universal hash family is a probability space of functions, which have uniform distribution and form k -wise independent random variables.

They can often be used in place of classic (or cryptographic) hash functions and allow the rigorous analysis of the performance of randomized algorithms and data structures that rely on hash functions.

In 1981 Wegman and Carter [4] introduced a generic construction for such families with arbitrary k using polynomials over a finite field. This entry contains a formalization of them and establishes the property of k -universality.

To be useful the formalization also provides an explicit construction of finite fields using the factor ring of integers modulo a prime. Additionally, some generic results about independent families are shown that might be of independent interest.

1 Introduction and Definition

theory *Definitions*

imports *HOL-Probability.Independent-Family*
begin

Universal hash families are commonly used in randomized algorithms and data structures to randomize the input of algorithms, such that probabilistic methods can be employed without requiring any assumptions about the input distribution.

If we regard a family of hash functions from a domain D to a finite range R as a uniform probability space, then the family is k -universal if:

- For each $x \in D$ the evaluation of the functions at x forms a uniformly distributed random variable on R .
- The evaluation random variables for k or fewer distinct domain elements form an independent family of random variables.

This definition closely follows the definition from Vadhan [3, §3.5.5], with the minor modification that independence is required not only for exactly k , but also for *fewer* than k distinct domain elements. The correction is due to the fact that in the corner case where D has fewer than k elements, the second part of their definition becomes void. In the formalization this helps avoid an unnecessary assumption in the theorems.

The following definition introduces the notion of k -wise independent random variables:

definition (in *prob-space*) *k-wise-indep-vars* **where**
k-wise-indep-vars k $M' X I =$
 $(\forall J \subseteq I. \text{card } J \leq k \longrightarrow \text{finite } J \longrightarrow \text{indep-vars } M' X J)$

lemma (in *prob-space*) *k-wise-indep-vars-subset*:
assumes *k-wise-indep-vars* k $M' X I$
assumes $J \subseteq I$
assumes *finite* J
assumes $\text{card } J \leq k$
shows *indep-vars* $M' X J$
 $\langle \text{proof} \rangle$

Similarly for a finite non-empty set A the predicate *uniform-on* $X A$ indicates that the random variable is uniformly distributed on A :

definition (in *prob-space*) *uniform-on* $X A =$ (
 $\text{distr } M (\text{count-space } UNIV) X = \text{uniform-measure } (\text{count-space } UNIV) A \wedge$
 $A \neq \{\}$) \wedge *finite* $A \wedge$ *random-variable* (*count-space* $UNIV$) X)

lemma (in *prob-space*) *uniform-onD*:
assumes *uniform-on* $X A$
shows $\text{prob } \{\omega \in \text{space } M. X \omega \in B\} = \text{card } (A \cap B) / \text{card } A$
 $\langle \text{proof} \rangle$

With the two previous definitions it is possible to define the k -universality condition for a family of hash functions from D to R :

definition (in *prob-space*) *k-universal* $k X D R =$ (
k-wise-indep-vars k $(\lambda-. \text{count-space } UNIV) X D \wedge$
 $(\forall i \in D. \text{uniform-on } (X i) R)$)

Note: The definition is slightly more generic than the informal specification from above. This is because usually a family is formed by a single function with a variable seed parameter. Instead of choosing a random function from a probability space, a random seed is chosen from the probability space which parameterizes the hash function.

The following section contains some preliminary results about independent families of random variables. Section 3 introduces the Carter-Wegman hash family, which is an explicit construction of k -universal families for arbitrary k using polynomials over finite fields. The last section contains a proof that

the factor ring of the integers modulo a prime ideal is a finite field, followed by an isomorphic construction of prime fields over an initial segment of the natural numbers.

end

2 Preliminary Results

theory *Preliminary-Results*

imports

Definitions

HOL-Probability.Stream-Space

HOL-Probability.Probability-Mass-Function

begin

lemma *set-comp-image-cong*:

assumes $\bigwedge x. P\ x \implies f\ x = h\ (g\ x)$

shows $\{f\ x \mid x. P\ x\} = h\ \text{' } \{g\ x \mid x. P\ x\}$

<proof>

lemma (**in** *prob-space*) *k-wise-indep-vars-compose*:

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

assumes $\bigwedge i. i \in I \implies Y\ i \in \text{measurable}\ (M'\ i)\ (N\ i)$

shows *k-wise-indep-vars* $k\ N\ (\lambda i\ x. Y\ i\ (X\ i\ x))\ I$

<proof>

The following two lemmas are of independent interest, they help infer independence of events and random variables on distributions. (Candidates for *HOL-Probability.Independent-Family*).

lemma (**in** *prob-space*) *indep-sets-distr*:

fixes A

assumes *random-variable* $N\ f$

defines $F \equiv (\lambda i. (\lambda a. f\ \text{' } a \cap \text{space}\ M)\ \text{' } A\ i)$

assumes *indep-F*: *indep-sets* $F\ I$

assumes *sets-A*: $\bigwedge i. i \in I \implies A\ i \subseteq \text{sets}\ N$

shows *prob-space.indep-sets* $(\text{distr}\ M\ N\ f)\ A\ I$

<proof>

lemma (**in** *prob-space*) *indep-vars-distr*:

assumes $f \in \text{measurable}\ M\ N$

assumes $\bigwedge i. i \in I \implies X'\ i \in \text{measurable}\ N\ (M'\ i)$

assumes *indep-vars* $M'\ (\lambda i. (X'\ i) \circ f)\ I$

shows *prob-space.indep-vars* $(\text{distr}\ M\ N\ f)\ M'\ X'\ I$

<proof>

lemma *range-inter*: $\text{range}\ ((\cap)\ F) = \text{Pow}\ F$

<proof>

The singletons and the empty set form an intersection stable generator of a

countable discrete σ -algebra:

lemma *sigma-sets-singletons-and-empty*:

assumes *countable* M

shows *sigma-sets* M (*insert* $\{\}$ $((\lambda k. \{k\}) \text{ ` } M)) = \text{Pow } M$

<proof>

In some of the following theorems, the premise $M = \text{measure-pmf } p$ is used. This allows stating theorems that hold for pmfs more concisely, for example, instead of $\text{measure-pmf.prob } p \ A \leq \text{measure-pmf.prob } p \ B$ we can just write $M = \text{measure-pmf } p \implies \text{prob } A \leq \text{prob } B$ in the locale *prob-space*.

lemma *prob-space-restrict-space*:

assumes $[\text{simp}]: M = \text{measure-pmf } p$

shows *prob-space* (*restrict-space* M (*set-pmf* p))

<proof>

The abbreviation below is used to specify the discrete σ -algebra on *UNIV* as a measure space. It is used in places where the existing definitions, such as *indep-vars*, expect a measure space even though only a *measurable* space is really needed, i.e., in cases where the property is invariant with respect to the actual measure.

hide-const (**open**) *discrete*

abbreviation *discrete* $\equiv \text{count-space } \text{UNIV}$

lemma (**in** *prob-space*) *indep-vars-restrict-space*:

assumes $[\text{simp}]: M = \text{measure-pmf } p$

assumes

prob-space.indep-vars (*restrict-space* M (*set-pmf* p)) ($\lambda\cdot$. *discrete*) $X \ I$

shows *indep-vars* ($\lambda\cdot$. *discrete*) $X \ I$

<proof>

lemma (**in** *prob-space*) *measure-pmf-eq*:

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

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

shows $\text{prob } P = \text{prob } Q$

<proof>

The following lemma is an intro rule for the independence of random variables defined on pmfs. In that case it is possible, to check the independence of random variables point-wise.

The proof relies on the fact that the support of a pmf is countable and the σ -algebra of such a set can be generated by singletons.

lemma (**in** *prob-space*) *indep-vars-pmf*:

assumes $[\text{simp}]: M = \text{measure-pmf } p$

assumes $\bigwedge a \ J. J \subseteq I \implies \text{finite } J \implies$

$\text{prob } \{\omega. \forall i \in J. X \ i \ \omega = a \ i\} = (\prod i \in J. \text{prob } \{\omega. X \ i \ \omega = a \ i\})$

shows *indep-vars* ($\lambda\cdot$. *discrete*) $X \ I$

<proof>

lemma (in *prob-space*) *split-indep-events*:

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

assumes *indep-vars* ($\lambda i. \text{discrete}$) $X' I$

assumes $K \subseteq I$ *finite* K

shows $\text{prob } \{\omega. \forall x \in K. P x (X' x \omega)\} = (\prod x \in K. \text{prob } \{\omega. P x (X' x \omega)\})$

<proof>

lemma *pmf-of-set-eq-uniform*:

assumes *finite* A $A \neq \{\}$

shows $\text{measure-pmf } (\text{pmf-of-set } A) = \text{uniform-measure discrete } A$

<proof>

lemma (in *prob-space*) *uniform-onI*:

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

assumes *finite* A $A \neq \{\}$

assumes $\bigwedge a. \text{prob } \{\omega. X \omega = a\} = \text{indicator } A a / \text{card } A$

shows *uniform-on* $X A$

<proof>

end

3 Carter-Wegman Hash Family

theory *Carter-Wegman-Hash-Family*

imports

*Interpolation-Polynomials-HOL-Algebra.Interpolation-Polynomial-Cardinalities
Preliminary-Results*

begin

The Carter-Wegman hash family is a generic method to obtain k -universal hash families for arbitrary k . (There are faster solutions, such as tabulation hashing, which are limited to a specific k . See for example [2].)

The construction was described by Wegman and Carter [4], it is a hash family between the elements of a finite field and works by choosing randomly a polynomial over the field with degree less than k . The hash function is the evaluation of a such a polynomial.

Using the property that the fraction of polynomials interpolating a given set of $s \leq k$ points is $1 / \text{real } (\text{card } (\text{carrier } R))^s$, which is shown in [1], it is possible to obtain both that the hash functions are k -wise independent and uniformly distributed.

In the following two locales are introduced, the main reason for both is to make the statements of the theorems and proofs more concise. The first locale *poly-hash-family* fixes a finite ring R and the probability space of the polynomials of degree less than k . Because the ring is not a field, the family

is not yet k -universal, but it is still possible to state a few results such as the fact that the range of the hash function is a subset of the carrier of the ring.

The second locale *carter-wegman-hash-family* is an extension of the former with the assumption that R is a field with which the k -universality follows. The reason for using two separate locales is to support use cases, where the ring is only probably a field. For example if it is the set of integers modulo an approximate prime, in such a situation a subset of the properties of an algorithm using approximate primes would need to be verified even if R is only a ring.

definition (in *ring*) $hash\ x\ \omega = eval\ \omega\ x$

locale *poly-hash-family* = *ring* +
fixes $k :: nat$
assumes *finite-carrier[simp]*: *finite* (*carrier* R)
assumes *k-ge-0*: $k > 0$
begin

definition *space* **where** *space* = *bounded-degree-polynomials* $R\ k$

definition M **where** $M = measure-pmf$ (*pmf-of-set* *space*)

lemma *finite-space[simp]*: *finite* *space*
 $\langle proof \rangle$

lemma *non-empty-bounded-degree-polynomials[simp]*: *space* $\neq \{\}$
 $\langle proof \rangle$

This is to add *carrier-not-empty* to the simp set in the context of *poly-hash-family*:

lemma *non-empty-carrier[simp]*: *carrier* $R \neq \{\}$
 $\langle proof \rangle$

sublocale *prob-space* M
 $\langle proof \rangle$

lemma *hash-range[simp]*:
assumes $\omega \in space$
assumes $x \in carrier\ R$
shows $hash\ x\ \omega \in carrier\ R$
 $\langle proof \rangle$

lemma *hash-range-2*:
assumes $\omega \in space$
shows $(\lambda x. hash\ x\ \omega) \text{ 'carrier } R \subseteq carrier\ R$
 $\langle proof \rangle$

lemma *integrable-M[simp]*:
fixes $f :: 'a\ list \Rightarrow 'c :: \{banach, second-countable-topology\}$

shows *integrable* $M f$
 ⟨*proof*⟩

end

locale *carter-wegman-hash-family* = *poly-hash-family* +
assumes *field-R*: *field* R
begin
sublocale *field*
 ⟨*proof*⟩

abbreviation *field-size* \equiv *card* (*carrier* R)

lemma *poly-cards*:
assumes $K \subseteq$ *carrier* R
assumes *card* $K \leq k$
assumes $y \in K \subseteq$ (*carrier* R)
shows
card $\{\omega \in$ *space*. $(\forall k \in K. \text{eval } \omega \ k = y \ k)\}$ = *field-size* $^{(k - \text{card } K)}$
 ⟨*proof*⟩

lemma *poly-cards-single*:
assumes $x \in$ *carrier* R
assumes $y \in$ *carrier* R
shows *card* $\{\omega \in$ *space*. $\text{eval } \omega \ x = y\}$ = *field-size* $^{(k - 1)}$
 ⟨*proof*⟩

lemma *hash-prob*:
assumes $K \subseteq$ *carrier* R
assumes *card* $K \leq k$
assumes $y \in K \subseteq$ *carrier* R
shows
prob $\{\omega. (\forall x \in K. \text{hash } x \ \omega = y \ x)\}$ = $1 / (\text{real } \text{field-size})^{\text{card } K}$
 ⟨*proof*⟩

lemma *prob-single*:
assumes $x \in$ *carrier* R $y \in$ *carrier* R
shows *prob* $\{\omega. \text{hash } x \ \omega = y\}$ = $1 / (\text{real } \text{field-size})$
 ⟨*proof*⟩

lemma *prob-range*:
assumes [*simp*]: $x \in$ *carrier* R
shows *prob* $\{\omega. \text{hash } x \ \omega \in A\}$ = *card* ($A \cap$ *carrier* R) / *field-size*
 ⟨*proof*⟩

lemma *indep*:
assumes $J \subseteq$ *carrier* R
assumes *card* $J \leq k$
shows *indep-vars* (λ -. *discrete*) *hash* J

<proof>

lemma *k-wise-indep*:

k-wise-indep-vars k (λ -. *discrete*) *hash* (*carrier* R)

<proof>

lemma *inj-if-degree-1*:

assumes $\omega \in$ *space*

assumes *degree* $\omega = 1$

shows *inj-on* (λx . *hash* $x \omega$) (*carrier* R)

<proof>

lemma *uniform*:

assumes $i \in$ *carrier* R

shows *uniform-on* (*hash* i) (*carrier* R)

<proof>

This the main result of this section - the Carter-Wegman hash family is k -universal.

theorem *k-universal*:

k-universal k *hash* (*carrier* R) (*carrier* R)

<proof>

end

lemma *poly-hash-familyI*:

assumes *ring* R

assumes *finite* (*carrier* R)

assumes $0 < k$

shows *poly-hash-family* $R k$

<proof>

lemma *carter-wegman-hash-familyI*:

assumes *field* F

assumes *finite* (*carrier* F)

assumes $0 < k$

shows *carter-wegman-hash-family* $F k$

<proof>

lemma *hash-k-wise-indep*:

assumes *field* $F \wedge$ *finite* (*carrier* F)

assumes $1 \leq n$

shows

prob-space.k-wise-indep-vars (*pmf-of-set* (*bounded-degree-polynomials* $F n$)) n

(λ -. *pmf-of-set* (*carrier* F)) (*ring.hash* F) (*carrier* F)

<proof>

lemma *hash-prob-single*:

assumes *field* $F \wedge$ *finite* (*carrier* F)

```

assumes  $x \in \text{carrier } F$ 
assumes  $1 \leq n$ 
assumes  $y \in \text{carrier } F$ 
shows
   $\mathcal{P}(\omega \text{ in pmf-of-set (bounded-degree-polynomials } F \ n). \text{ ring.hash } F \ x \ \omega = y)$ 
   $= 1 / (\text{real (card (carrier } F)))$ 
<proof>

end

```

4 Finite Fields

```

theory Field
imports
  Finite-Fields.Ring-Characteristic
  HOL-Algebra.Ring-Divisibility
  HOL-Algebra.IntRing
begin

```

In some applications it is more convenient to work with natural numbers instead of $ZFact\ p$ whose elements are cosets. To support that use case the following definition introduces an additive and multiplicative structure on $\{..

\}$. After verifying that the function *zfact-iso* and its inverse are homomorphisms, the ring and field property can be transferred from $ZFact\ p$ to the structure on $\{..

\}$.

```

lemma zfact-iso-0:
assumes  $n > 0$ 
shows  $\text{zfact-iso } n \ 0 = \mathbf{0}_{ZFact\ (int\ n)}$ 
<proof>

```

```

lemma zfact-prime-is-field:
assumes Factorial-Ring.prime ( $p :: \text{nat}$ )
shows field ( $ZFact\ (int\ p)$ )
<proof>

```

```

definition mod-ring ::  $\text{nat} \Rightarrow \text{nat ring}$ 
where mod-ring  $n = \langle$ 
  carrier =  $\{..

\}$ ,
  mult =  $(\lambda\ x\ y. (x * y) \text{ mod } n)$ ,
  one =  $1$ ,
  zero =  $0$ ,
  add =  $(\lambda\ x\ y. (x + y) \text{ mod } n)$   $\rangle$ 

```

```

definition zfact-iso-inv ::  $\text{nat} \Rightarrow \text{int set} \Rightarrow \text{nat}$  where
  zfact-iso-inv  $p = \text{inv-into } \{..

\} (\text{zfact-iso } p)$ 

```

```

lemma zfact-iso-inv-0:
assumes n-ge-0:  $n > 0$ 

```

shows $zfact\text{-}iso\text{-}inv\ n\ \mathbf{0}\ ZFact\ (int\ n) = 0$
 $\langle proof \rangle$

lemma *zfact-coset*:
assumes $n\text{-}ge\text{-}0: n > 0$
assumes $x \in carrier\ (ZFact\ (int\ n))$
defines $I \equiv Idl_{\mathbb{Z}}\ \{int\ n\}$
shows $x = I\ +>_{\mathbb{Z}}\ (int\ (zfact\text{-}iso\text{-}inv\ n\ x))$
 $\langle proof \rangle$

lemma *zfact-iso-inv-is-ring-iso*:
assumes $n\text{-}ge\text{-}1: n > 1$
shows $zfact\text{-}iso\text{-}inv\ n \in ring\text{-}iso\ (ZFact\ (int\ n))\ (mod\text{-}ring\ n)$
 $\langle proof \rangle$

lemma *mod-ring-finite*:
 $finite\ (carrier\ (mod\text{-}ring\ n))$
 $\langle proof \rangle$

lemma *mod-ring-carr*:
 $x \in carrier\ (mod\text{-}ring\ n) \longleftrightarrow x < n$
 $\langle proof \rangle$

lemma *mod-ring-is-crng*:
assumes $n\text{-}ge\text{-}1: n > 1$
shows $crng\ (mod\text{-}ring\ n)$
 $\langle proof \rangle$

lemma *zfact-iso-is-ring-iso*:
assumes $n\text{-}ge\text{-}1: n > 1$
shows $zfact\text{-}iso\ n \in ring\text{-}iso\ (mod\text{-}ring\ n)\ (ZFact\ (int\ n))$
 $\langle proof \rangle$

If p is a prime than *mod-ring* p is a field:

lemma *mod-ring-is-field*:
assumes *Factorial-Ring.prime* p
shows $field\ (mod\text{-}ring\ p)$
 $\langle proof \rangle$

end

References

- [1] E. Karayel. Interpolation polynomials (in hol-algebra). *Archive of Formal Proofs*, Jan. 2022. https://isa-afp.org/entries/Interpolation_Polynomials_HOL_Algebra.html, Formal proof development.

- [2] M. Thorup and Y. Zhang. Tabulation based 5-universal hashing and linear probing. In *Proceedings of the Meeting on Algorithm Engineering & Experiments, ALENEX '10*, pages 62–76, USA, 2010. Society for Industrial and Applied Mathematics.
- [3] S. P. Vadhan. Pseudorandomness. *Foundations and Trends® in Theoretical Computer Science*, 7(1-3):1–336, 2012.
- [4] M. N. Wegman and J. L. Carter. New hash functions and their use in authentication and set equality. *Journal of Computer and System Sciences*, 22(3):265–279, 1981.