

3.3 VARIANCE AND STANDARD DEVIATION RECAP

ANNA KARLIN
MOST SLIDES BY ALEX TSUN

AGENDA

- VARIANCE
- INDEPENDENCE OF RANDOM VARIABLES
- PROPERTIES OF VARIANCE

VARIANCE AND STANDARD DEVIATION (SD)

Variance: The variance of a random variable X is

MORE USEFUL

$$\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$$

The variance is always nonnegative since we take an expectation of a nonnegative random variable $(X - E[X])^2$. We can also show that for any scalars $a, b \in \mathbb{R}$,

$$\text{Var}(aX + b) = a^2 \text{Var}(X)$$

Standard Deviation (SD): The standard deviation of a random variable X is

$$\sigma_X = \sqrt{\text{Var}(X)}$$

We want this because the units of variance are squared in terms of the original variable X , and this “undo’s” our squaring, returning the units to the same as X .

RANDOM VARIABLES AND INDEPENDENCE

Random variable X and event E are independent if the event E is independent of the event $\{X=x\}$ (for any fixed x), i.e.

$$\forall x \ P(X = x \text{ and } E) = P(X=x) \cdot P(E)$$

Two random variables X and Y are independent if the events $\{X=x\}$ and $\{Y=y\}$ are independent for any fixed x, y , i.e.

$$\forall x, y \ P(X = x \text{ and } Y=y) = P(X=x) \cdot P(Y=y)$$

Intuition as before: knowing X doesn't help you guess Y or E and vice versa.

INDEPENDENT VS DEPENDENT R.V.S

- Dependent r.v.s can reinforce/cancel/correlate in arbitrary ways.
- Independent r.v.s are, well, independent.

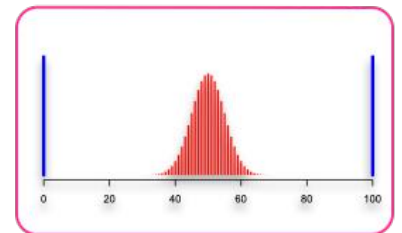
Example:

$$Z = X_1 + X_2 + \dots + X_n$$

X_i is indicator r.v. with probability 1/2 of being 1.

versus

$$W = n X_1$$



IMPORTANT FACTS ABOUT INDEPENDENT RANDOM VARIABLES

Theorem: If X & Y are independent, then $E[X \cdot Y] = E[X] \cdot E[Y]$

Theorem: If X and Y are independent, then

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

Corollary: If $X_1 + X_2 + \dots + X_n$ are mutually independent then

$$\text{Var}[X_1 + X_2 + \dots + X_n] = \text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n]$$

E[XY] FOR INDEPENDENT RANDOM VARIABLES

- Theorem: If X & Y are independent, then $E[X \cdot Y] = E[X] \cdot E[Y]$
- Proof:

Let $x_i, y_i, i = 1, 2, \dots$ be the possible values of X, Y .

$$\begin{aligned} E[X \cdot Y] &= \sum_i \sum_j x_i \cdot y_j \cdot P(X = x_i \wedge Y = y_j) \quad \leftarrow \text{independence} \\ &= \sum_i \sum_j x_i \cdot y_j \cdot P(X = x_i) \cdot P(Y = y_j) \\ &= \sum_i x_i \cdot P(X = x_i) \cdot \left(\sum_j y_j \cdot P(Y = y_j) \right) \\ &= E[X] \cdot E[Y] \end{aligned}$$

Note: *NOT* true in general; see earlier example $E[X^2] \neq E[X]^2$

VARIANCE OF A SUM OF INDEPENDENT R.V.S

Theorem: If X and Y are independent, then

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

Proof:

$$\begin{aligned}\text{Var}[X + Y] &= E[(X + Y)^2] - (E[X + Y])^2 \\ &= E[X^2 + 2XY + Y^2] - (E[X] + E[Y])^2 \\ &= E[X^2] + 2E[XY] + E[Y^2] - ((E[X])^2 + 2E[X]E[Y] + (E[Y])^2) \\ &= E[X^2] - (E[X])^2 + E[Y^2] - (E[Y])^2 + 2(E[XY] - E[X]E[Y]) \\ &= \text{Var}[X] + \text{Var}[Y] + 2(\overbrace{E[X]E[Y]}^{\leftarrow}) - E[X]E[Y] \\ &= \text{Var}[X] + \text{Var}[Y]\end{aligned}$$

x



BLOOM FILTERS

ANNA KARLIN

MOST SLIDES BY SHREYA JAYARAMAN, LUXI WANG, ALEX TSUN

HASHING

BASIC PROBLEM

Problem: Store a subset S of a large set U .

Example. U = set of 128 bit strings
 S = subset of strings of interest

$$|U| \approx 2^{128}$$

$$|S| \approx 1000$$

Two goals:

1. **Constant-time** answering of queries “Is $x \in S$?”
2. **Minimize storage** requirements.

NAÏVE SOLUTION – CONSTANT TIME

Idea: Represent S as an array A with 2^{128} entries.

$$A[x] = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases}$$

$S = \{0, 2, \dots, K\}$



0	1	2	...	K	...		
1	0	1	0	1	...	0	0

Membership test: To check $x \in S$ just check whether $A[x] = 1$.

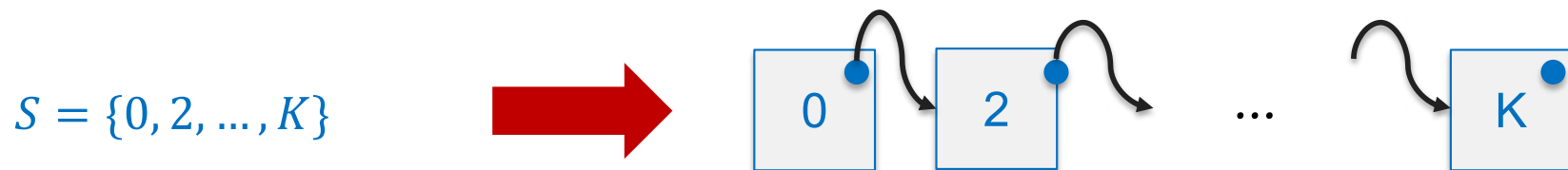
→ constant time! 👍 😊

Storage: Require storing 2^{64} bits, even for small S .



NAÏVE SOLUTION – SMALL STORAGE

Idea: Represent S as a list with $|S|$ entries.



Storage: Grows with $|S|$ only



Membership test: Check $x \in S$ requires time linear in $|S|$

(Can be made logarithmic by using a tree)

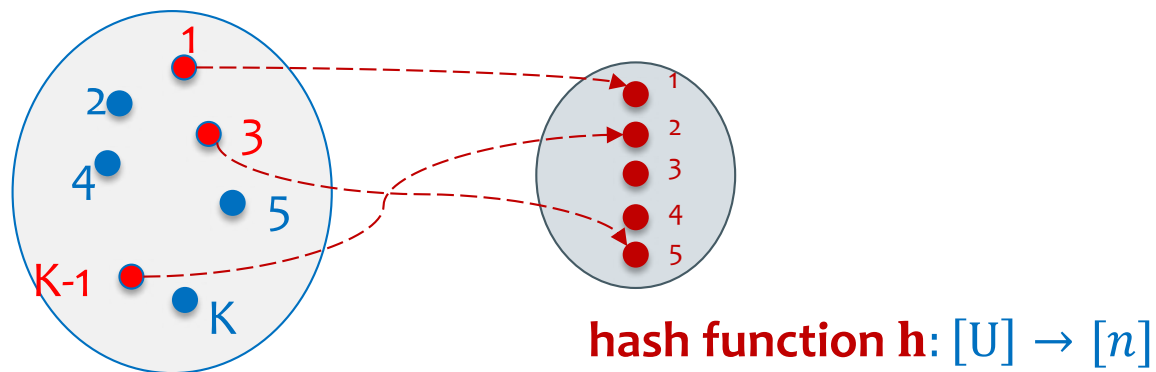


HASH TABLE

Idea: Map elements in S into an array A using a hash function

Membership test: To check $x \in S$
just check whether $A[\mathbf{h}(x)] = x$

Storage: n elements

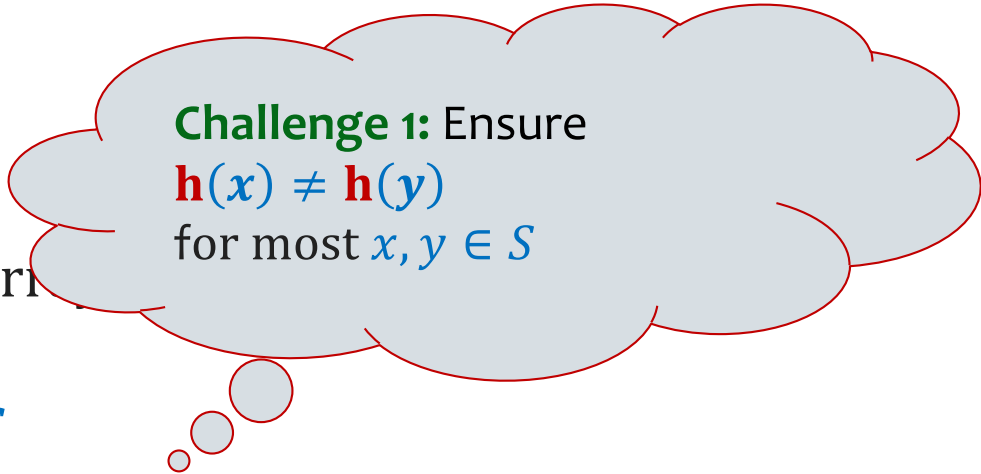


HASH TABLE

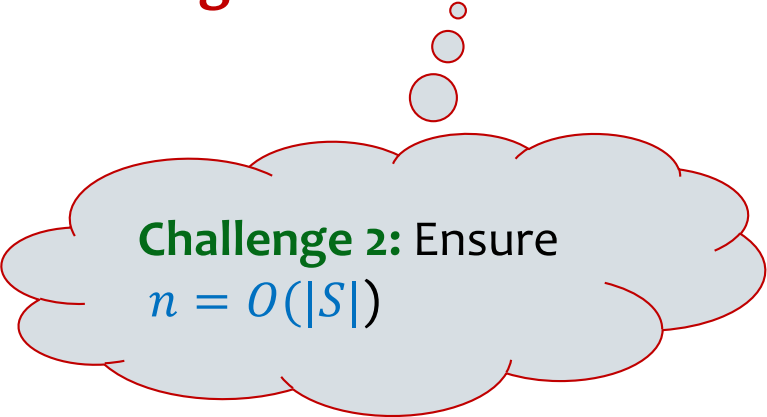
Idea: Map elements in S into an array

Membership test: To check $x \in S$
just check whether $A[h(x)] = x$

Storage: n elements



Challenge 1: Ensure
 $h(x) \neq h(y)$
for most $x, y \in S$



Challenge 2: Ensure
 $n = O(|S|)$

HASHING - COLLISIONS

- **Collisions** occur when two elements of set map to the same location in the hash table.
- Common solution: chaining - at each location (bucket) in the table, keep linked list of all elements that hash there.
- Want: hash function that distributes the elements of S well across hash table locations. Ideally uniform distribution!

SUMMARY



Hash Tables

- They store the data itself
- With a good hash function, the data is well distributed in the table and lookup times are small.
- However, they need at least as much space as all the data being stored
- E.g. storing strings, or IP addresses or long DNA sequences.

BLOOM FILTERS: MOTIVATION



- Large universe of possible data items.
- Data items are large (say 128 bits or more)
- Hash table is stored on disk or across network, so any lookup is expensive.
- Many (if not nearly all) of the lookups return “Not found”.

Altogether, this is bad. You're wasting **a lot of time and space** doing lookups for items that aren't even present.

BLOOM FILTERS: MOTIVATION



- Large universe of possible data items.
- Hash table is stored on disk or in network, so any lookup is expensive.
- Many (if not most) of the lookups return “Not found”.

Altogether, this is bad. You’re wasting **a lot of time and space** doing lookups for items that aren’t even present.

Examples:

- Google Chrome: wants to warn you if you’re trying to access a malicious URL. Keep hash table of malicious URLs.
- Network routers: want to track source IP addresses of certain packets, .e.g., blocked IP addresses.

BLOOM FILTERS: MOTIVATION



- Probabilistic data structure.
- Close cousins of hash tables.
- Ridiculously space efficient
- To get that, make occasional errors, specifically false positives.

Typical implementation: only 8 bits per element!

BLOOM FILTERS

BLOOM FILTERS



- Stores information about a set of elements.
- Supports two operations:
 1. **add(x)** - adds x to bloom filter
 2. **contains(x)** - returns true if x in bloom filter, otherwise returns false
 - a. If return false, **definitely** not in bloom filter.
 - b. If return true, **possibly** in the structure (some false positives).

BLOOM FILTERS



- Why accept false positives?
 - **Speed** – both operations very very fast.
 - **Space** – requires a miniscule amount of space relative to storing all the actual items that have been added.
- Often just 8 bits per inserted item!

BLOOM FILTERS: INITIALIZATION

Number of
hash
functions

Size of array
associated to
each hash
function.

```
function INITIALIZE(k,m)
  for  $i = 1, \dots, k$ : do
     $t_i =$  new bit vector of  $m$  0's
```

for each hash
function,
initialize an
empty bit vector
of size m

BLOOM FILTERS: EXAMPLE

bloom filter t with $m = 5$ that uses $k = 3$ hash functions

```
function INITIALIZE(k,m)
  for  $i = 1, \dots, k$ : do
     $t_i =$  new bit vector of  $m$  0's
```

Index →	0	1	2	3	4
t_1	0	0	0	0	0
t_2	0	0	0	0	0
t_3	0	0	0	0	0

BLOOM FILTERS: ADD

```
function ADD(x)  
  for  $i = 1, \dots, k$ : do  
     $t_i[h_i(x)] = 1$ 
```

for each hash
function h_i

$h_i(x) \rightarrow$ result of hash function h_i on x

BLOOM FILTERS: ADD

```
function ADD(x)  
  for  $i = 1, \dots, k$  do  
     $t_i[h_i(x)] = 1$ 
```

for each hash
function h_i

Index into i th bit-vector, at index
produced by hash function and set to 1

BLOOM FILTERS: EXAMPLE

bloom filter t with $m = 5$ that uses $k = 3$ hash functions

`add("thisisavirus.com")`

```
function ADD( $x$ )  
  for  $i = 1, \dots, k$ : do  
     $t_i[h_i(x)] = 1$ 
```

Index →	0	1	2	3	4
t_1	0	0	0	0	0
t_2	0	0	0	0	0
t_3	0	0	0	0	0

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function ADD(x)
  for  $i = 1, \dots, k$ : do
     $t_i[h_i(x)] = 1$ 
```

add("thisisavirus.com")

$h_1(\text{"thisisavirus.com"}) \rightarrow 2$

Index →	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	0	0	0	0
t_3	0	0	0	0	0

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function ADD( $x$ )  
  for  $i = 1, \dots, k$ : do  
     $t_i[h_i(x)] = 1$ 
```

add("thisisavirus.com")

$h_1(\text{"thisisavirus.com"}) \rightarrow 2$

$h_2(\text{"thisisavirus.com"}) \rightarrow 1$

Index \rightarrow	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	0

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function ADD(x)
  for  $i = 1, \dots, k$ : do
     $t_i[h_i(x)] = 1$ 
```

add("thisisavirus.com")

$h_1(\text{"thisisavirus.com"}) \rightarrow 2$

$h_2(\text{"thisisavirus.com"}) \rightarrow 1$

$h_3(\text{"thisisavirus.com"}) \rightarrow 4$

Index →	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: CONTAINS

```
function CONTAINS(x)
```

```
return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

Returns True if the bit vector for each hash function has bit 1 at index determined by that hash function, otherwise returns False

BLOOM FILTERS: EXAMPLE

bloom filter t with $m = 5$ that uses $k = 3$ hash functions

`contains("thisisavirus.com")`

```
function CONTAINS(x)  
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

Index →	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function CONTAINS(x)
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

contains("thisisavirus.com")

$h_1(\text{"thisisavirus.com"}) \rightarrow 2$

Index →	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function CONTAINS(x)
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

contains("thisisavirus.com")

h_1 ("thisisavirus.com") \rightarrow 2

h_2 ("thisisavirus.com") \rightarrow 1

Index \rightarrow	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function CONTAINS(x)
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

True

contains("thisisavirus.com")

h_1 ("thisisavirus.com") \rightarrow 2

h_2 ("thisisavirus.com") \rightarrow 1

h_3 ("thisisavirus.com") \rightarrow 4

Index \rightarrow	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

contains("thisisavirus.com")

```
function CONTAINS(x)
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

True

h_1 ("thisisavirus.com") \rightarrow 2

h_2 ("thisisavirus.com") \rightarrow 1

h_3 ("thisisavirus.com") \rightarrow 4

Since all conditions satisfied, returns True (correctly)

Index \rightarrow	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: FALSE POSITIVES

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

`add("totallynotsuspicious.com")`

function `ADD(x)`

for $i = 1, \dots, k$: **do**

$t_i[h_i(x)] = 1$

Index →	0	1	2	3	4
t_1	0	0	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: FALSE POSITIVES

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

`add("totallynotsuspicious.com")`

`h1("totallynotsuspicious.com") → 1`

function `ADD(x)`

for $i = 1, \dots, k$: **do**

$t_i[h_i(x)] = 1$

Index →	0	1	2	3	4
t_1	0	1	1	0	0
t_2	0	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: FALSE POSITIVES

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function ADD( $x$ )  
  for  $i = 1, \dots, k$ : do  
     $t_i[h_i(x)] = 1$ 
```

```
add("totallynotsuspicious.com")
```

```
 $h_1$ ("totallynotsuspicious.com")  $\rightarrow$  1
```

```
 $h_2$ ("totallynotsuspicious.com")  $\rightarrow$  0
```

Index \rightarrow	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: FALSE POSITIVES

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function ADD( $x$ )  
  for  $i = 1, \dots, k$ : do  
     $t_i[h_i(x)] = 1$ 
```

```
add("totallynotsuspicious.com")
```

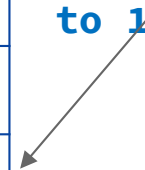
```
 $h_1$ ("totallynotsuspicious.com")  $\rightarrow$  1
```

```
 $h_2$ ("totallynotsuspicious.com")  $\rightarrow$  0
```

```
 $h_3$ ("totallynotsuspicious.com")  $\rightarrow$  4
```

Index \rightarrow	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

Collision, is
already set
to 1



BLOOM FILTERS: FALSE POSITIVES

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function ADD( $x$ )  
  for  $i = 1, \dots, k$ : do  
     $t_i[h_i(x)] = 1$ 
```

```
add("totallynotsuspicious.com")
```

```
 $h_1$ ("totallynotsuspicious.com")  $\rightarrow$  1
```

```
 $h_2$ ("totallynotsuspicious.com")  $\rightarrow$  0
```

```
 $h_3$ ("totallynotsuspicious.com")  $\rightarrow$  4
```

Index \rightarrow	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

contains("verynormalsite.com")

function CONTAINS(x)

return $t_1[h_1(x)] \wedge t_2[h_2(x)] \wedge \dots \wedge t_k[h_k(x)]$

Index →	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

contains("verynormalsite.com")

$h_1(\text{"verynormalsite.com"}) \rightarrow 2$

function CONTAINS(x)

return $t_1[h_1(x)] \wedge t_2[h_2(x)] \wedge \dots \wedge t_k[h_k(x)]$

True

Index →	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function CONTAINS(x)  
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

contains("verynormalsite.com")

$h_1(\text{"verynormalsite.com"}) \rightarrow 2$

$h_2(\text{"verynormalsite.com"}) \rightarrow 0$

Index →	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function CONTAINS(x)  
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

True

contains("verynormalsite.com")

$h_1(\text{"verynormalsite.com"}) \rightarrow 2$

$h_2(\text{"verynormalsite.com"}) \rightarrow 0$

$h_3(\text{"verynormalsite.com"}) \rightarrow 4$

Index →	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: EXAMPLE

bloom filter t of length $m = 5$ that uses $k = 3$ hash functions

```
function CONTAINS(x)  
  return  $t_1[h_1(x)] == 1 \wedge t_2[h_2(x)] == 1 \wedge \dots \wedge t_k[h_k(x)] == 1$ 
```

True

True

True

contains("verynormalsite.com")

$h_1(\text{"verynormalsite.com"}) \rightarrow 2$

$h_2(\text{"verynormalsite.com"}) \rightarrow 0$

$h_3(\text{"verynormalsite.com"}) \rightarrow 4$

Since all conditions satisfied, returns True (incorrectly)

Index →	0	1	2	3	4
t_1	0	1	1	0	0
t_2	1	1	0	0	0
t_3	0	0	0	0	1

BLOOM FILTERS: SUMMARY



- An empty bloom filter is an empty $k \times m$ bit array with all values initialized to zeros
 - k = number of hash functions
 - m = size of each array in the bloom filter
- `add(x)` runs in $O(k)$ time
- `contains(x)` runs in $O(k)$ time
- requires $O(km)$ space (in bits!)
- Probability of false positives from collisions can be reduced by increasing the size of the bloom filter

BLOOM FILTERS: APPLICATION



- Google Chrome has a database of malicious URLs, but it takes a long time to query.
- Want an in-browser structure, so needs to be efficient and be space-efficient
- Want it so that can check if a URL is in structure:
 - If return False, then definitely not in the structure (don't need to do expensive database lookup, website is safe)
 - If return True, the URL may or may not be in the structure. Have to perform expensive lookup in this rare case.

FALSE POSITIVE PROBABILITY

COMPARISON WITH HASH TABLES - SPACE



- Google storing 5 million URLs, each URL 40 bytes.
- Bloom filter with $k=8$ and $m = 10,000,000$.

Hash Table

Bloom Filter

COMPARISON WITH HASH TABLES - TIME



- Say avg user visits 100,000 URLs in a year, of which 2,000 are malicious.
- 0.5 seconds to do lookup in the database, 1ms for lookup in Bloom filter.
- Suppose the false positive rate is 2%

Hash Table

Bloom Filter

BLOOM FILTERS: MANY APPLICATIONS



- Any scenario where space and efficiency are important.
- Used a lot in networking
- In distributed systems when want to check consistency of data across different locations, might send a Bloom filter rather than the full set of data being stored.
- Google BigTable uses Bloom filters to reduce the disk lookups for non-existent rows and columns
- Internet routers often use Bloom filters to track blocked IP addresses.
- And on and on...

BLOOM FILTERS TYPICAL EXAMPLE...



of randomized algorithms and randomized data structures.

- Simple
 - Fast
 - Efficient
 - Elegant
 - Useful!
-
- You'll be implementing Bloom filters on pset 4. Enjoy!