# COMPSCI 514: Algorithms for Data Science

Cameron Musco

University of Massachusetts Amherst. Spring 2026.

Lecture 6

## Logistics

- Problem Set 1 solutions have been posted.
- Problem Set 2 will be posted in the next couple of days and due before the first midterm.

## Quiz Poll

Lecture Pacing:

- Way too Fast – 8
- A Bit too Fast – 15
- Just Right – 23
- Too Slow – 1
- Way Too Slow – 0

Last Class: $M = \max |X_i|$

- Exponential concentration bounds — Sums of $n$ bounded independent r.v.s. (not necessarily i.i.d)
- Bernstein inequality and the Chernoff bound
- Connection to the central limit theorem.
- Application to random hashing – $O(\log n)$ maximum load per bucket when hashing $n$ items into $n$ buckets.

This Class:

- Bloom filters: random hashing to maintain a large set in small space.

4

## Approximately Maintaining a Set

Want to store a set $S$ of items from a massive universe of possible items (e.g., images, text documents, IP addresses).

## Approximately Maintaining a Set

Want to store a set *S* of items from a massive universe of possible items (e.g., images, text documents, IP addresses).

**Goal:** support *insert*(*x*) to add *x* to the set and *query*(*x*) to check if *x* is in the set. Both in *O*(1) time.

~~*delete*(*x*)~~

## Approximately Maintaining a Set

Want to store a set *S* of items from a massive universe of possible items (e.g., images, text documents, IP addresses).

**Goal:** support *insert*(*x*) to add *x* to the set and *query*(*x*) to check if *x* is in the set. Both in *O*(1) time. What data structure(s) solves this problem?

- hash table
- binary tree, linked list, array

## Approximately Maintaining a Set

Want to store a set *S* of items from a massive universe of possible items (e.g., images, text documents, IP addresses).

**Goal:** support *insert*(*x*) to add *x* to the set and *query*(*x*) to check if *x* is in the set. Both in $O(1)$ time. What data structure solves this problem?

- Allow small probability $\delta > 0$ of false positives. I.e., for any *x*,

$$\delta = .01$$

$$\Pr(query(x) = 1 \text{ and } x \notin S) \leq \delta.$$

## Approximately Maintaining a Set

Want to store a set *S* of items from a massive universe of possible items (e.g., images, text documents, IP addresses).

**Goal:** support *insert*(*x*) to add *x* to the set and *query*(*x*) to check if *x* is in the set. Both in *O*(1) time. What data structure solves this problem?

- Allow small probability $\delta > 0$ of false positives. I.e., for any *x*,

$$\Pr(query(x) = 1 \text{ and } x \notin S) \leq \delta.$$

**Solution:** Bloom filters (repeated random hashing). Will use much less space than a hash table.
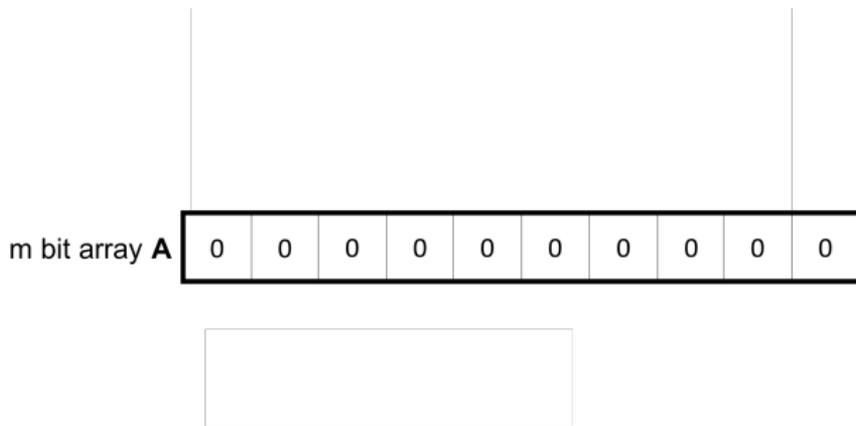
## Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.
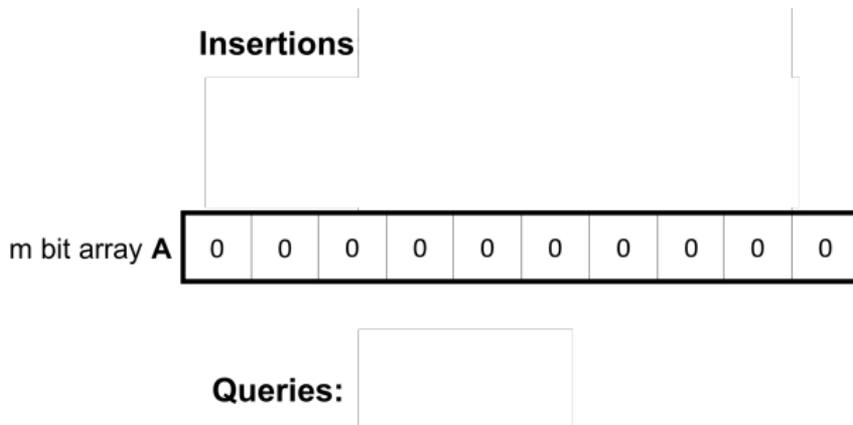
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

## Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

| m bit array **A** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|

## Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.
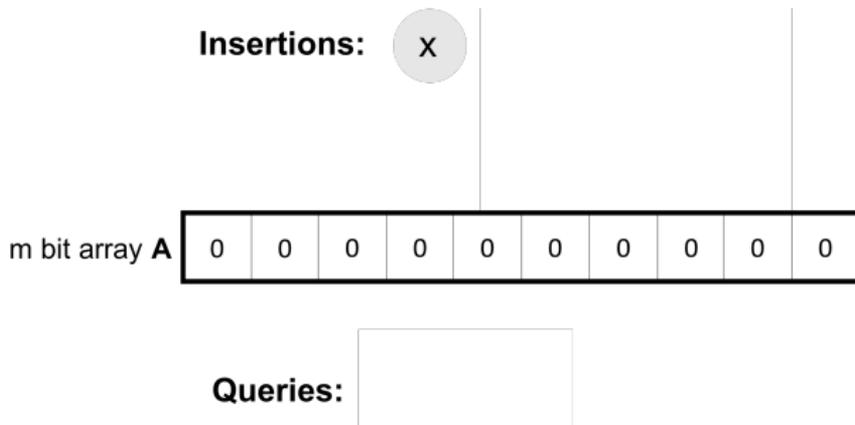
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

**Insertions**

m bit array **A**

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

**Queries:**

## Bloom Filters

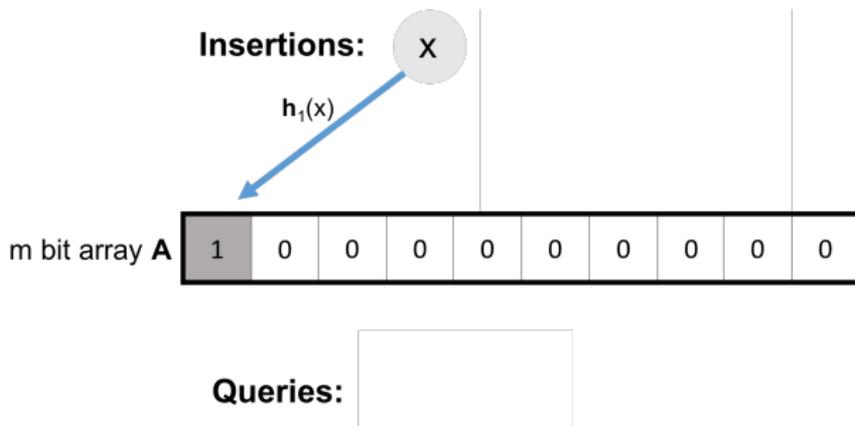Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

**Insertions:** x

m bit array **A**

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

**Queries:**

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \rightarrow [m]$.
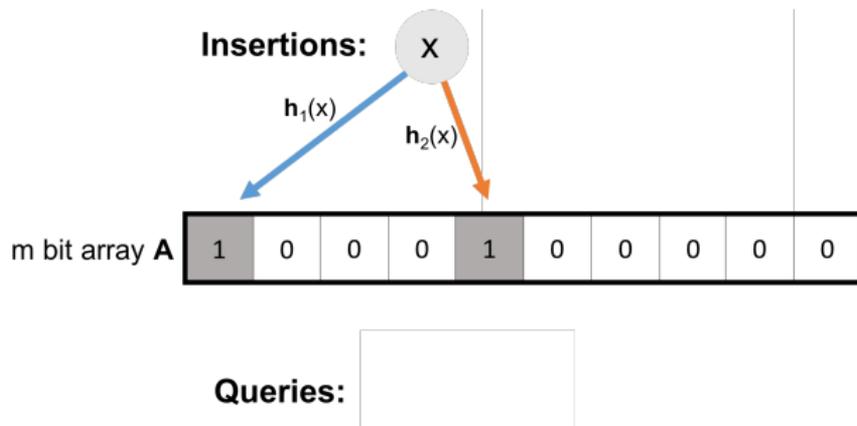
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
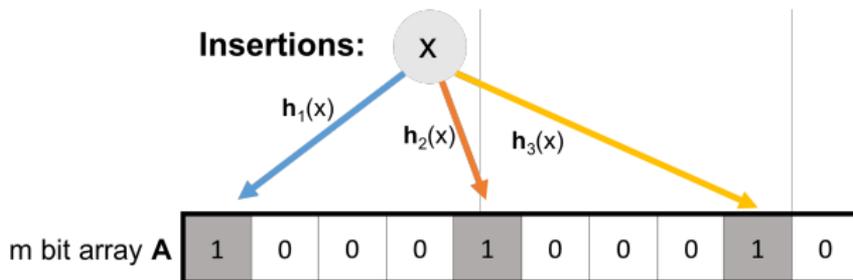- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

**Insertions:** X

$h_1(x)$

m bit array **A** | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Queries:**

# Bloom Filters

Chose $k$ independent random hash functions $\mathbf{h}_1, \ldots, \mathbf{h}_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] := 1$.
- *query(x)*: return 1 only if $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \rightarrow [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
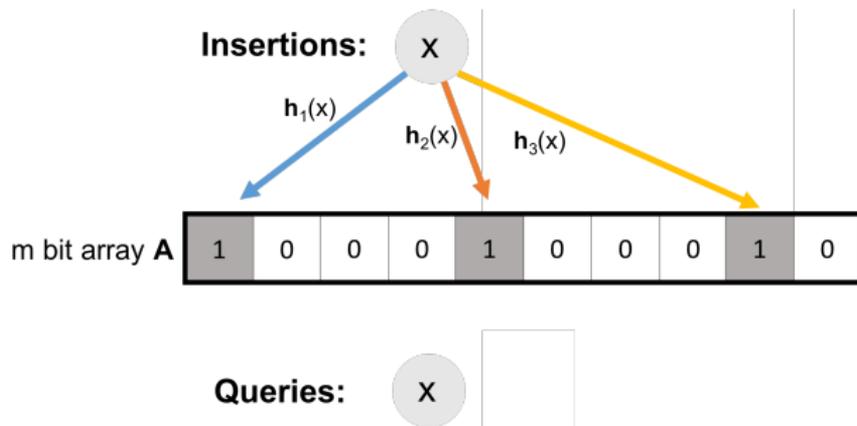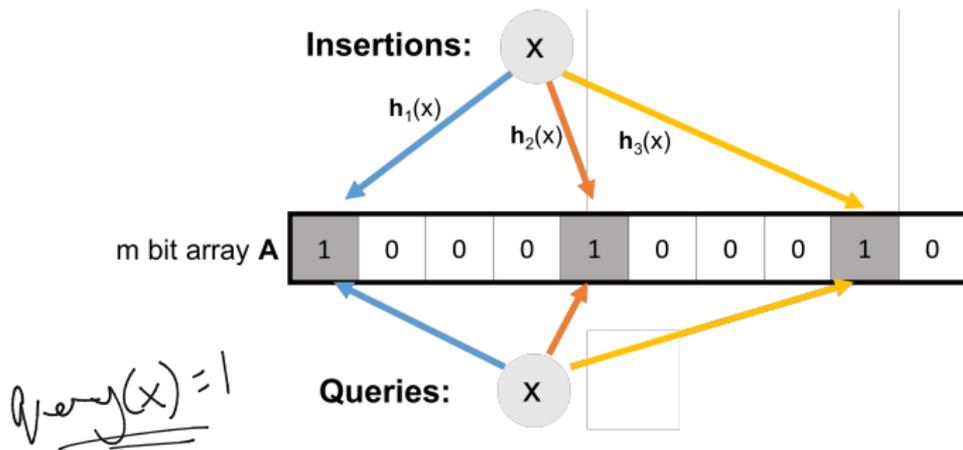- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

$K = 3$

# Bloom Filters

Chose $k$ independent random hash functions $\mathbf{h}_1, \ldots, \mathbf{h}_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] := 1$.
- *query(x)*: return 1 only if $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $\mathbf{h}_1, \ldots, \mathbf{h}_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] := 1$.
- *query(x)*: return 1 only if $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.
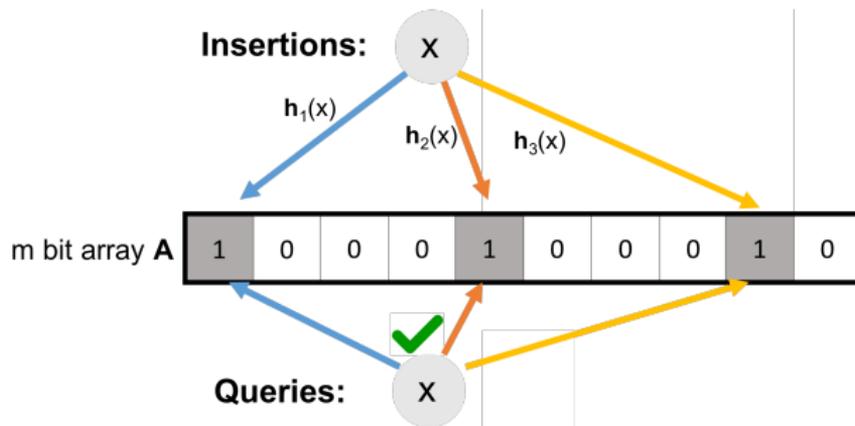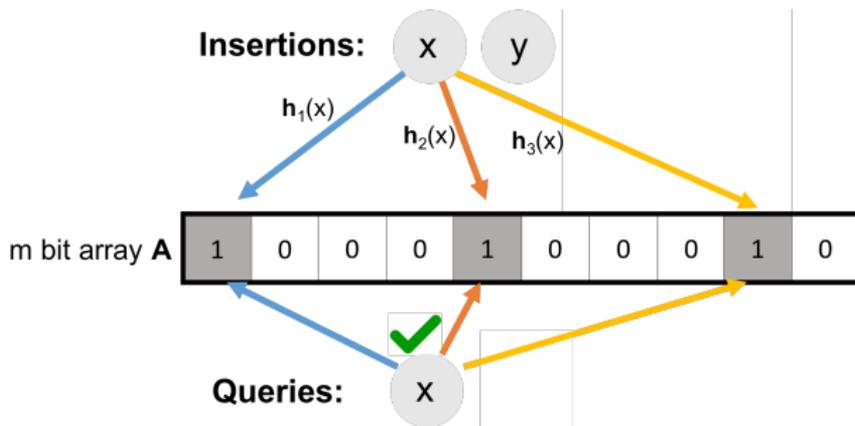
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $\mathbf{h}_1, \ldots, \mathbf{h}_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] := 1$.
- *query(x)*: return 1 only if $A[\mathbf{h}_1(x)] = \ldots = A[\mathbf{h}_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.
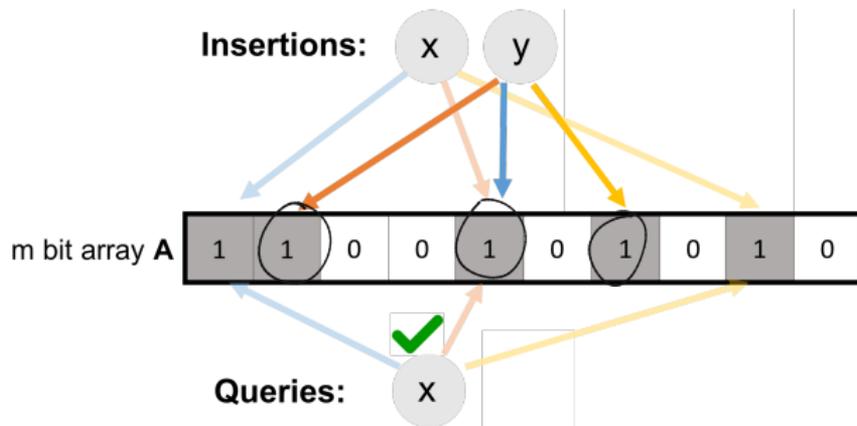
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \rightarrow [m]$.
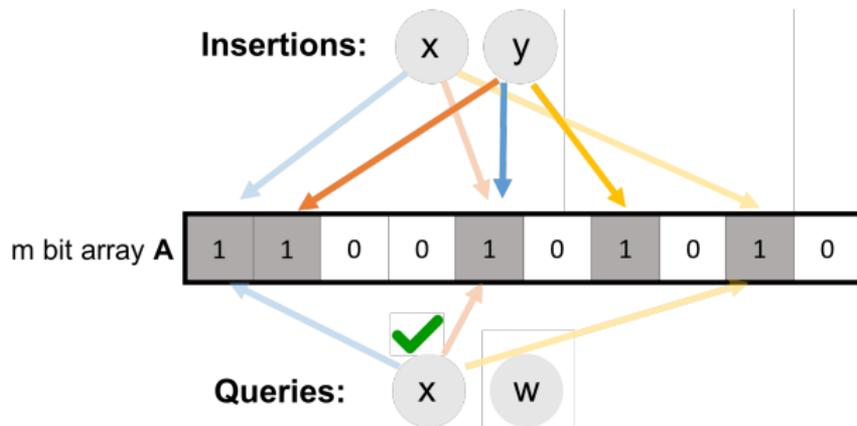
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.
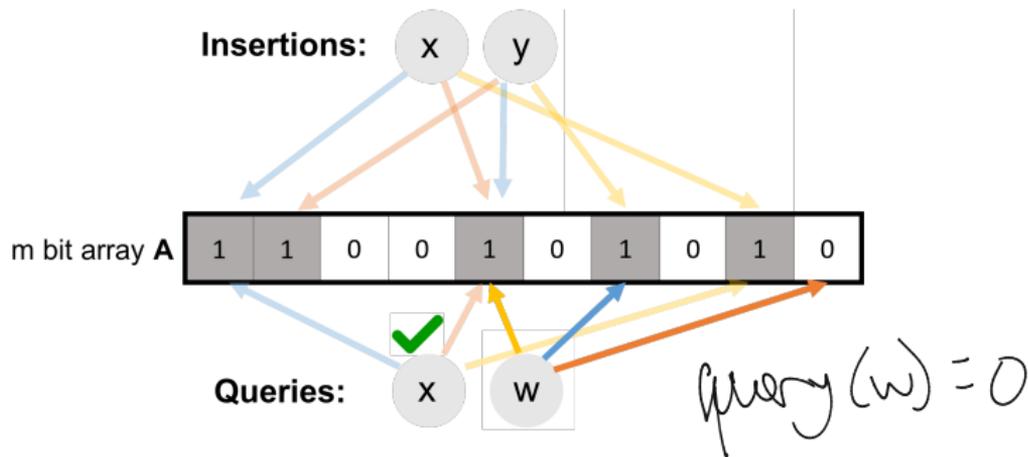
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.



Insertions: x   y

m bit array **A**  | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |

Queries: x   w

$query(w) = 0$

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \to [m]$.

- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.



**Insertions:** x y

m bit array **A**: 1 1 0 0 1 0 1 0 1 0

**Queries:** x ✅ w ❌

# Bloom Filters

Chose $k$ independent random hash functions $h_1, \ldots, h_k$ mapping the universe of elements $U \rightarrow [m]$.
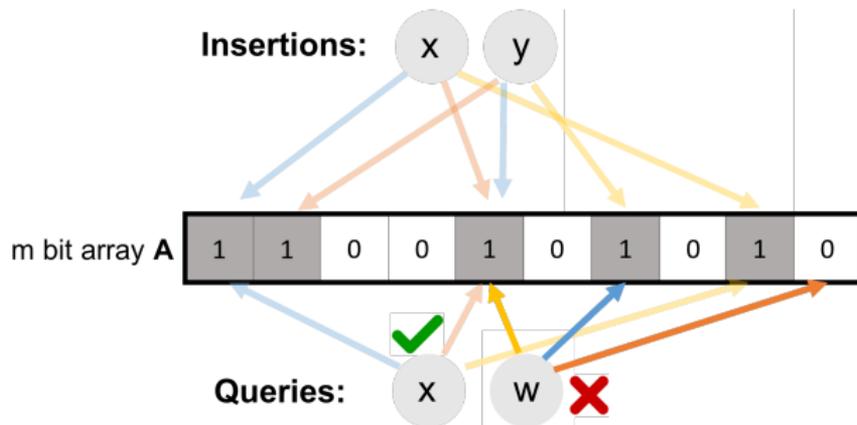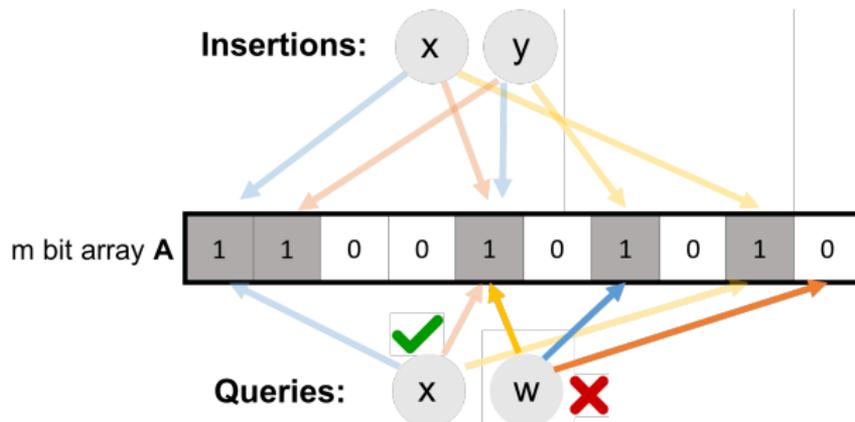
- Maintain an array $A$ containing $m$ bits, all initially 0.
- *insert(x)*: set all bits $A[h_1(x)] = \ldots = A[h_k(x)] := 1$.
- *query(x)*: return 1 only if $A[h_1(x)] = \ldots = A[h_k(x)] = 1$.



No false negatives. False positives more likely with more insertions.

Akamai (Boston-based company serving 15 − 30% of all web traffic) applies bloom filters to prevent caching of 'one-hit-wonders' – pages only visited once fill over 75% of cache.

Akamai (Boston-based company serving $15 - 30\%$ of all web traffic) applies bloom filters to prevent caching of 'one-hit-wonders' – pages only visited once fill over 75% of cache.



- When url $x$ comes in, if $query(x) = 1$, cache the page at $x$. If not, run $insert(x)$ so that if it comes in again, it will be cached.
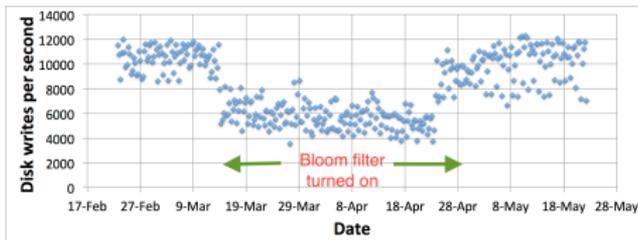
# Applications: Caching

Akamai (Boston-based company serving $15 - 30\%$ of all web traffic) applies bloom filters to prevent caching of 'one-hit-wonders' – pages only visited once fill over 75% of cache.



- When url $x$ comes in, if $query(x) = 1$, cache the page at $x$. If not, run $insert(x)$ so that if it comes in again, it will be cached.

- **False positive:** A new url (possible one-hit-wonder) is cached. If the bloom filter has a false positive rate of $\delta = .05$, the number of cached one-hit-wonders will be reduced by at least 95%.

## Applications: Databases

Distributed database systems, including Google Bigtable, Apache HBase, Apache Cassandra, and PostgreSQL use bloom filters to prevent expensive lookups of non-existent data.

## Applications: Databases

Distributed database systems, including Google Bigtable, Apache HBase, Apache Cassandra, and PostgreSQL use bloom filters to prevent expensive lookups of non-existent data.

## Applications: Databases

Distributed database systems, including Google Bigtable, Apache HBase, Apache Cassandra, and PostgreSQL use bloom filters to prevent expensive lookups of non-existent data.

Movies



Users

- When a new rating is inserted for ($user_x$, $movie_y$), add ($user_x$, $movie_y$) to a bloom filter.
- Before reading ($user_x$, $movie_y$) (possibly via an out of memory access), check the bloom filter, which is stored in memory.

## Applications: Databases

Distributed database systems, including Google Bigtable, Apache HBase, Apache Cassandra, and PostgreSQL use bloom filters to prevent expensive lookups of non-existent data.



- When a new rating is inserted for ($user_x$, $movie_y$), add ($user_x$, $movie_y$) to a bloom filter.
- Before reading ($user_x$, $movie_y$) (possibly via an out of memory access), check the bloom filter, which is stored in memory.
- **False positive:** A read is made to a possibly empty cell. A $\delta = .05$ false positive rate gives a 95% reduction in these empty reads.

## More Applications

- **Database Joins:** Quickly eliminate most keys in one column that don't correspond to keys in another.
- **Recommendation systems:** Bloom filters are used to prevent showing users the same recommendations twice.
- **Spam/Fraud Detection**:
    - Bit.ly and Google Chrome use bloom filters to quickly check if a url maps to a flagged site and prevent a user from following it.
    - Can be used to detect repeat clicks on the same ad from a single IP-address, which may be the result of fraud.
- **Digital Currency:** Some Bitcoin clients use bloom filters to quickly pare down the full transaction log to transactions involving bitcoin addresses that are relevant to them (SPV: simplified payment verification).

## Analysis

For a bloom filter with $m$ bits and $k$ hash functions, the insertion and query time is $O(k)$.

For a bloom filter with $m$ bits and $k$ hash functions, the insertion and query time is $O(k)$. How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

$$m \rightarrow \infty$$
$$k \rightarrow \infty$$

$$FPR \rightarrow 0$$

$$\left(1 - \frac{1}{m}\right)^k \gtrsim 1 - \frac{k}{m}$$

For a bloom filter with $m$ bits and $k$ hash functions, the insertion and query time is $O(k)$. How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?



$n = 1$

$k = 1$

0?

$X_1 \cdots = X_n$

(is this bit 0?)

$n = 1, k = 1$

$$\frac{m-1}{m} = 1 - \frac{1}{m}$$

$n = 1 \quad k \geqslant 1$

$$\left(1 - \frac{1}{m}\right)^k$$

$$1 \times \frac{k}{m}$$

$n \geqslant 1, \quad k \geqslant 1$

$$\left(1 - \frac{1}{m}\right)^{kn}$$

10

## Analysis

For a bloom filter with $m$ bits and $k$ hash functions, the insertion and query time is $O(k)$. How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0? $n \times k$ total hashes must not hit bit $i$.

$$\Pr(A[i] = 0) = \Pr\left(h_1(x_1) \neq i \cap \ldots \cap h_k(x_1) \neq i \right.$$
$$\left. \cap \underbrace{h_1(x_2) \neq i \ldots \cap h_k(x_2) \neq i}_{k \cdot n} \cap \ldots \right)$$

For a bloom filter with $m$ bits and $k$ hash functions, the insertion and query time is $O(k)$. How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0? $n \times k$ total hashes must not hit bit $i$.

$$
\begin{aligned}
\Pr(A[i] = 0) &= \Pr\left(\mathsf{h}_1(x_1) \neq i \cap \ldots \cap \mathsf{h}_k(x_k) \neq i \right. \\
&\qquad \left. \cap \mathsf{h}_1(x_2) \neq i \ldots \cap \mathsf{h}_k(x_2) \neq i \cap \ldots \right) \\
&= \underbrace{\Pr\left(\mathsf{h}_1(x_1) \neq i\right) \times \ldots \times \Pr\left(\mathsf{h}_k(x_1) \neq i\right) \times \Pr\left(\mathsf{h}_1(x_2) \neq i\right) \ldots}_{k \cdot n \text{ events each occuring with probability } 1-1/m}
\end{aligned}
$$

For a bloom filter with $m$ bits and $k$ hash functions, the insertion and query time is $O(k)$. How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0? $n \times k$ total hashes must not hit bit $i$.

$$\Pr(A[i] = 0) = \Pr\left( h_1(x_1) \neq i \cap \ldots \cap h_k(x_k) \neq i \right.$$
$$\left. \cap h_1(x_2) \neq i \ldots \cap h_k(x_2) \neq i \cap \ldots \right)$$
$$= \underbrace{\Pr\left( h_1(x_1) \neq i \right) \times \ldots \times \Pr\left( h_k(x_1) \neq i \right) \times \Pr\left( h_1(x_2) \neq i \right) \ldots}_{k \cdot n \text{ events each occuring with probability } 1 - 1/m}$$
$$= \left( 1 - \frac{1}{m} \right)^{kn}$$

*(handwritten annotations):*

$k = 3$

$[ \ ,^{\nearrow} \ ,^{\nwarrow} \ ]$

$\left( 1 - \frac{m-k}{m} \right)^n$

$\left( \frac{k}{m} \right)^n \longrightarrow$ true if $\Pr\left( \text{inserting 1 item does not } A[i]:1 \right) = \frac{k}{m} \quad \left( 1 - \frac{1}{m} \right)^k$

10

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn}$$

$n \to \infty \quad \Pr \to 1$

$k \to \infty \quad \Pr \to 0$

$n \to \infty \quad \Pr \to 0$

---

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $h_1, \ldots h_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\left(1 - \frac{1}{100}\right)^{100} \approx .37$$
$$\approx \frac{1}{e}$$

$$m=2$$
$$\left(1 - \frac{1}{2}\right)^2 = \frac{1}{4} \neq \frac{1}{e}$$

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

$$\lim_{m \to \infty} \left(1 - \frac{1}{m}\right)^m = \frac{1}{e}$$

$$\left(1 - \frac{1}{m}\right)^{kn} = \left[\left(1 - \frac{1}{m}\right)^m\right]^{\frac{kn}{m}} \approx \left(\frac{1}{e}\right)^{\frac{kn}{m}} = e^{-kn/m}$$

---

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $\mathbf{h}_1, \ldots \mathbf{h}_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

---

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $\mathbf{h}_1, \ldots \mathbf{h}_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

# Analysis

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

$$\Pr\left(A[\mathbf{h}_1(w)] = \ldots = A[\mathbf{h}_k(w)] = 1\right)$$
$$= \Pr(A[\mathbf{h}_1(w)] = 1) \times \ldots \times \Pr(A[\mathbf{h}_k(w)] = 1)$$

$$FPR = \left(1 - e^{-kn/m}\right)^k$$

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $\mathbf{h}_1, \ldots \mathbf{h}_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

$$
\begin{aligned}
\Pr\left(A[h_1(w)] = \ldots = A[h_k(w)] = 1\right) \\
= \Pr(A[h_1(w)] = 1) \times \ldots \times \Pr(A[h_k(w)] = 1) \\
= \left(1 - e^{-\frac{kn}{m}}\right)^k
\end{aligned}
$$

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $h_1, \ldots h_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

$$
\begin{aligned}
\Pr\left(A[\mathbf{h}_1(w)] = \ldots = A[\mathbf{h}_k(w)] = 1\right) \\
= \Pr(A[\mathbf{h}_1(w)] = 1) \times \ldots \times \Pr(A[\mathbf{h}_k(w)] = 1) \\
= \left(1 - e^{-\frac{kn}{m}}\right)^k \quad \textcolor{red}{\textbf{Actually Incorrect!}}
\end{aligned}
$$

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $\mathbf{h}_1, \ldots \mathbf{h}_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

How does the false positive rate $\delta$ depend on $m$, $k$, and the number of items inserted?

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

$$\begin{aligned}
\Pr\left(A[h_1(w)] = \ldots = A[h_k(w)] = 1\right) & \\
= \Pr(A[h_1(w)] = 1) \times \ldots \times \Pr(A[h_k(w)] = 1) & \\
= \left(1 - e^{-\frac{kn}{m}}\right)^k \quad \text{\textbf{Actually Incorrect!}} \text{ Dependent events.}
\end{aligned}$$

$n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $h_1, \ldots h_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

## Correct Analysis Sketch

Step 1: To avoid dependence issues, condition on the event that the $A$ has $t$ zeros in it after $n$ insertions, for some $t \leq m$. For a non-inserted element $w$, after conditioning on this event we correctly have:

$$\Pr(A[h_1(w)] = \ldots = A[h_k(w)] = 1)$$
$$= \Pr(A[h_1(w)] = 1) \times \ldots \times \Pr(A[h_k(w)] = 1).$$

I.e., the events $A[h_1(w)] = 1$,..., $A[h_k(w)] = 1$ are independent conditioned on the number of bits set in $A$. Why?

- Conditioned on this event, for any $j$, since $h_j$ is a fully random hash function, $\Pr(A[h_j(w)] = 1) = 1 - \frac{t}{m}$.
- Thus conditioned on this event, the false positive rate is $\left(1 - \frac{t}{m}\right)^k$.
- It remains to show that $\frac{t}{m} \approx e^{-\frac{kn}{m}}$ with high probability. We already have that $\mathbb{E}[\frac{t}{m}] = \frac{1}{m} \sum_{i=1}^m \Pr(A[i] = 0) \approx e^{-\frac{kn}{m}}$.

## Correct Analysis Sketch

Need to show that the number of zeros $t$ in $A$ after $n$ insertions is bounded by $O\left(e^{-\frac{kn}{m}}\right)$ with high probability.

Can apply Theorem 2 of:
http://cglab.ca/~morin/publications/ds/bloom-submitted.pdf

## Analysis

**Step 1**: What is the probability that after inserting $n$ elements, the $i^{th}$ bit of the array $A$ is still 0?

$$\Pr(A[i] = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-\frac{kn}{m}}$$

**Step 2**: What is the probability that querying a new item $w$ gives a false positive?

$$\Pr\left(A[h_1(w)] = \ldots = A[h_k(w)] = 1\right)$$
$$= \Pr(A[h_1(w)] = 1) \times \ldots \times \Pr(A[h_k(w)] = 1)$$
$$FPR = \left(1 - e^{-\frac{kn}{m}}\right)^k$$

> $n$: total number items in filter, $m$: number of bits in filter, $k$: number of random hash functions, $h_1, \ldots h_k$: hash functions, $A$: bit array, $\delta$: false positive rate.

## Optimizing Parameters

False Positive Rate: with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$.
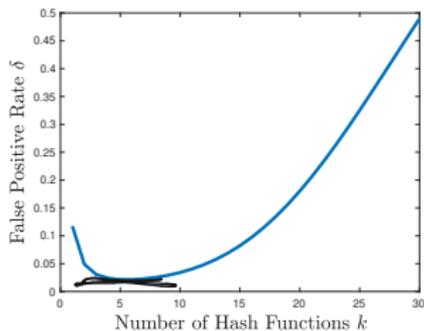
**False Positive Rate:** with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$. How should we set $k$ to minimize the FPR given a fixed amount of space $m$?

$k \uparrow$      filter fills with 1s   FPR $\uparrow$

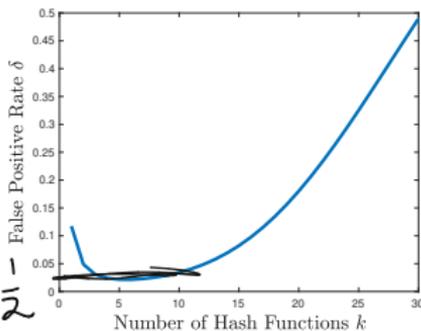$k \uparrow$   check more positions   FRR $\downarrow$

**False Positive Rate:** with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$. How should we set $k$ to minimize the FPR given a fixed amount of space $m$?

**False Positive Rate:** with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$ How should we set $k$ to minimize the FPR given a fixed amount of space $m$?

$$\delta = \underbrace{\left(1 - e^{\cdot \frac{kn}{m}}\right)}_{(\text{prob. entry} = 1)}^k$$

wur $k = \ln 2 \cdot \frac{m}{n}$

$$\text{Prob}(\text{entry} = 1) = \frac{1}{2}$$

$$k = \ln 2 \cdot \frac{m}{n}$$

$$FPR^* = \left(1 - e^{-\ln 2 \cdot \frac{m}{n} \cdot \frac{n}{m}}\right)^k$$

$$= \left(\frac{1}{2}\right)^k$$

at optimal Bloom Filter $\approx 1/2$ full



False Positive Rate $\delta$ vs. Number of Hash Functions $k$

- Can differentiate to show optimal number of hashes is $k = \ln 2 \cdot \frac{m}{n}$.
- Balances filling up the array vs. having enough hashes so that even when the array is pretty full, a new item is unlikely to yield a false positive.

15

## False Positive Rate

False Positive Rate: with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^{k}$.

Movies



Users

$10^9$ (user, movie) pairs

- Say we have 100 million users, each who have rated 10 movies.
- $n = 10^9 = n$ (user,movie) pairs with non-empty ratings.

## False Positive Rate

False Positive Rate: with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$.



Movies

Users

- Say we have 100 million users, each who have rated 10 movies.
- $n = 10^9 = n$ (user,movie) pairs with non-empty ratings.
- Allocate $m = 8n = 8 \times 10^9$ bits for a Bloom filter (1 GB).

## False Positive Rate

False Positive Rate: with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$.



Movies

Users

- Say we have 100 million users, each who have rated 10 movies.
- $n = 10^9 = n$ (user,movie) pairs with non-empty ratings.
- Allocate $m = 8n = 8 \times 10^9$ bits for a Bloom filter (1 GB).
- Set $k = \ln 2 \cdot \frac{m}{n} = 5.54 \approx 6$.

16

## False Positive Rate

False Positive Rate: with $m$ bits of storage, $k$ hash functions, and $n$ items inserted $\delta \approx \left(1 - e^{\frac{-kn}{m}}\right)^k$.



Movies

Users

- Say we have 100 million users, each who have rated 10 movies.
- $n = 10^9 = n$ (user,movie) pairs with non-empty ratings.
- Allocate $m = 8n = 8 \times 10^9$ bits for a Bloom filter (1 GB).
- Set $k = \ln 2 \cdot \frac{m}{n} = 5.54 \approx 6$.
- False positive rate is $\approx \left(1 - e^{-k \cdot \frac{n}{m}}\right)^k \approx \frac{1}{2^k} \approx \frac{1}{2^{5.54}} = .021$.
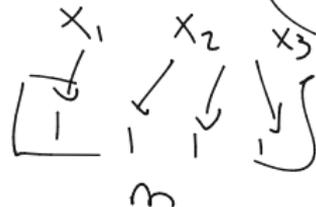
## Bloom Filter Note

An observation about Bloom filter space complexity:

$$\text{False Positive Rate: } \delta \approx \left(1 - e^{-\frac{kn}{m}}\right)^{k}.$$

For an $m$-bit bloom filter holding $n$ items, optimal number of hash functions $k$ is: $k = \ln 2 \cdot \frac{m}{n}$.

An observation about Bloom filter space complexity:

$$\text{False Positive Rate: } \delta \approx \left(1 - e^{-\frac{kn}{m}}\right)^k.$$

For an $m$-bit bloom filter holding $n$ items, optimal number of hash functions $k$ is: $k = \ln 2 \cdot \frac{m}{n}$.

**Think Pair Share:** If we want a false positive rate $< \frac{1}{2}$ how big does $m$ need to be in comparison to $n$?

$$m = O(\log n), \ m = O(\sqrt{n}), \ m = O(n), \ m = O(n^2)?$$

## Bloom Filter Note

An observation about Bloom filter space complexity:

$$\text{False Positive Rate: } \delta \approx \left(1 - e^{-\frac{kn}{m}}\right)^k.$$

For an $m$-bit bloom filter holding $n$ items, optimal number of hash functions $k$ is: $k = \ln 2 \cdot \frac{m}{n}$.

**Think Pair Share:** If we want a false positive rate $< \frac{1}{2}$ how big does $m$ need to be in comparison to $n$?

$$m = O(\log n), \ m = O(\sqrt{n}), \ m = O(n), \ m = O(n^2)?$$

If $m = \frac{n}{\ln 2}$, optimal $k = 1$, and failure rate is:

$$\delta = \left(1 - e^{-\frac{n/\ln 2}{n}}\right)^1 = \left(1 - \frac{1}{2}\right)^1 = \frac{1}{2}.$$

## Bloom Filter Note

An observation about Bloom filter space complexity:

$$\text{False Positive Rate: } \delta \approx \left(1 - e^{-\frac{kn}{m}}\right)^k.$$

For an $m$-bit bloom filter holding $n$ items, optimal number of hash functions $k$ is: $k = \ln 2 \cdot \frac{m}{n}$.

**Think Pair Share:** If we want a false positive rate $< \frac{1}{2}$ how big does $m$ need to be in comparison to $n$?

$$m = O(\log n), \ m = O(\sqrt{n}), \ m = O(n), \ m = O(n^2)?$$

If $m = \frac{n}{\ln 2}$, optimal $k = 1$, and failure rate is:

$$\delta = \left(1 - e^{-\frac{n/\ln 2}{n}}\right)^1 = \left(1 - \frac{1}{2}\right)^1 = \frac{1}{2}.$$

I.e., storing $n$ items in a bloom filter requires $O(n)$ space. So what's the point? Truly $O(n)$ bits, rather than $O(n \cdot \text{item size})$.

Questions on Bloom Filters?