

Bloom

Filter

Bloom Filter Operations

Bloom Filter Operations

$B.add(x)$

Bloom Filter Operations

$B.add(x)$

x in B ?

Bloom Filter Operations

$B.add(x)$

No deletion

$x \text{ in } B?$

Bloom Filter Operations

$B.add(x)$

No deletion

$x \text{ in } B?$

No key/value pairs

Bloom Filter Operations

$B.add(x)$

No deletion

$x \text{ in } B ?$

No key/value pairs

Main Points

Bloom Filter Operations

$B.add(x)$

No deletion

$x \text{ in } B ?$

No key/value pairs

Main Points

- Bloom filters are very small,
typically smaller than the data stored

Bloom Filter Operations

$B.add(x)$

No deletion

$x \text{ in } B?$

No key/value pairs

Main Points

- Bloom filters are very small, typically smaller than the data stored
- There is one-sided error, Bloom filters may say yes to $x \text{ in } B?$ when the answer is no.

Why is small important?

Why is small important?

- A key way to deal with big data is to make it smaller.

Why is small important?

- A key way to deal with big data is to make it smaller.
- Smaller may allow the data to be

Why is small important?

- A key way to deal with big data is to make it smaller.
- Smaller may allow the data to be
 - In cache rather than RAM

Why is small important?

- A key way to deal with big data is to make it smaller.
- Smaller may allow the data to be
 - In Cache rather than RAM
 - In RAM rather than Disc

Why is small important?

- A key way to deal with big data is to make it smaller.
- Smaller may allow the data to be
 - In cache rather than RAM
 - In RAM rather than Disc
 - On a local machine rather than a server

Why is small important?

- A key way to deal with big data is to make it smaller.
- Smaller may allow the data to be
 - In Cache rather than RAM
 - In RAM rather than Disc
 - On a local machine rather than a server
 - On a mobile Device

When is one-sided error OK?

When is One-sided error OK?

Example: A database of Bad URLs, when
a user types in one you want a warning
with perhaps some info

When is One-sided error OK?

Example: A database of Bad URLs, when a user types in one you went a warning with perhaps some info

Option 1:



Put the DB
on each user's
device

When is One-sided error OK?

Example: A database of Bad URLs, when a user types in one you want a warning with perhaps some info

Option 1:



Put the DB
on each user's
device

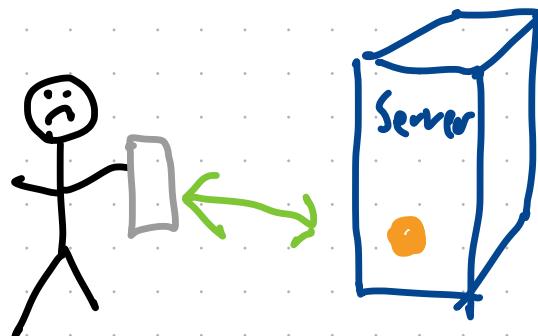
When is One-sided error OK?

Example: A database of Bad URLs, when a user types in one you want a warning with perhaps some info

Option 1:



Option 1:



Put the DB
on each user's
device

Put the DB
on a server

When is One-sided error OK?

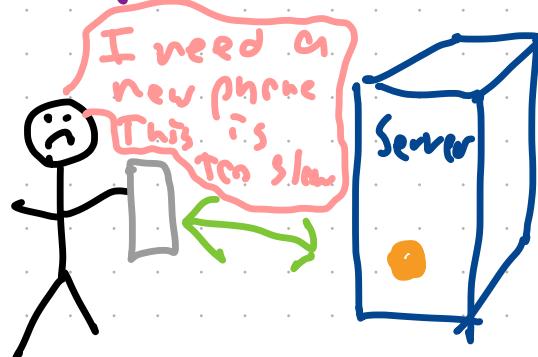
Example: A database of Bad URLs, when a user types in one you want a warning with perhaps some info

Option 1:



Put the DB
on each user's
device

Option 1:



Put the DB
on a server

When is One-sided error OK?

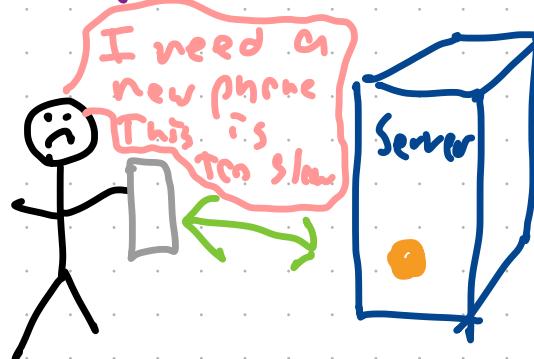
Example: A database of Bad URLs, when a user types in one you want a warning with perhaps some info

Option 1:



Put the DB
on each user's
device

Option 1:



Put the DB
on a server

Option 3:



Put a Bloom filter
on the phone and
a DB on a server.

Bloom Filters, Take 1

Bloom Filters, Take 1

- A Bloom filter is a Bit Array, call it A

Bloom Filters, Take 1

- A Bloom filter is a Bit Array, call it A
- Let K denote the size of the Bitarray
(usually $K = cn$, for a small c , e.g. 10)

Bloom Filters, Take 1

- A Bloom filter is a Bit Array, call it A
- Let K denote the size of the Bitarray
(usually $K = cn$, for a small c , e.g. 10)
- To insert x , set $A[\text{hash}(x) \% K] = 1$

Bloom Filters, Take 1

- A Bloom filter is a Bit Array, call it A
- Let K denote the size of the Bitarray
(usually $K = cn$, for a small c , e.g. 10)
- To insert x , set $A[\text{hash}(x) \circ, K] = 1$
- To check if x is in the filter
we ask is $A[\text{hash}(x) \circ, K] == 1$?

Bloom Filters, Take 1

- A Bloom filter is a Bit Array, call it A
- Let K denote the size of the Bitarray
(usually $K = cn$, for a small c , e.g. 10)
- To insert x , set $A[\text{hash}(x) \circ, K] = 1$
- To check if x is in the filter
we ask is $A[\text{hash}(x) \circ, K] == 1$?
Chance of a false positive is high!

Bloom Filters, Take 1

- A Bloom filter is a Bit Array, call it A
- Let K denote the size of the Bitarray
(usually $K = cn$, for a small c , e.g. 10)
- To insert x , set $A[\text{hash}(x) \circ, K] = 1$
- To check if x is in the filter
we ask is $A[\text{hash}(x) \circ, K] == 1$?
Chance of a false positive is high!
Can we improve?

Bloom Filters, For Real

- A Bloom filter is a Bit Array, call it A
- Let K denote the size of the Bitarray
(usually $K = cn$, for a small c , e.g. 10)
- Let h denote the number of hash functions
 - To insert x , set $A[\text{hash}_i(x) \text{ } g, K] = 1$
for all i
 - To check if x is in the filter
we ask is $A[\text{hash}_i(x) \text{ } g, K] == 1?$
for all i ?

La de

Analysis:

Analysis:

- What is the false positive rate
Given K, n, h ?

Analysis:

- What is the false positive rate
Given K, n, h ?
- How should we choose the
right number of hash functions

Analysis:

- What is the false positive rate
Given K, n, h ?
- How should we choose the right number of hash functions
- What is the tradeoff between error rate and K (size)?

False Positive Rate

False Positive Rate

- N items inserted, each sets h random Bits
to one

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$P[\text{single Bit zero}] =$

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$$P[\text{single Bit zero}] = \left(1 - \frac{1}{K}\right)^{hn}$$

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$$P[\text{single Bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \boxed{\left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}}}$$

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$$P[\text{single Bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \boxed{\left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}}} \approx e^{\frac{hn}{K}}$$

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$$P[\text{single Bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \boxed{\left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}}} \approx e^{\frac{hn}{K}}$$

$P[\text{h Bits one}]$
Random
(false positive)

False Positive Rate

- N items inserted, each sets h random Bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$$P[\text{single Bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \boxed{\left(1 - \frac{1}{K}\right)^K}^{\frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h Bits one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{\frac{hn}{K}}}$$

False Positive Rate

- N items inserted, each sets h random bits to one
- What is the chance that if we look at h random bits, they are all one?
- What is the chance that if we look at a single random bit it is zero

$$P[\text{single Bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \boxed{\left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h Bits one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{\frac{hn}{K}}}$$

This is minimized when $h = \frac{K}{n} \ln 2$

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[h \stackrel{\text{Random}}{\text{Bits}} \text{ one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{hn/K}}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[K \stackrel{\text{Random}}{\text{Bits}} \text{ one}] \approx$$

with optimal h
(false positive)

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h bits one} \text{ (false positive)}] \stackrel{\text{Random}}{\approx} \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{hn/K}}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[K \text{ bits one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h$$

with optimal h
(false positive)

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h bits one} \underset{\text{(false positive)}}{\text{Random}}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{hn/K}}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[\text{h bits one} \underset{\text{(false positive)}}{\text{Random}}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h = \left(1 - e^{\left(\frac{K}{n} \ln 2\right) \frac{n}{K}}\right)^{\frac{K}{n} \ln 2}$$

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h bits one} \underset{\substack{\text{Random} \\ (\text{false positive})}}{\approx} (1 - e^{\frac{hn}{K}})^h \underset{\approx e}{\approx} e^{-\frac{hn}{K}}]$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[\text{h bits one} \underset{\substack{\text{Random} \\ (\text{false positive})}}{\approx} (1 - e^{\frac{hn}{K}})^h = (1 - e^{(\frac{K}{n} \ln 2) \frac{n}{K}})^{\frac{K}{n} \ln 2} = \left(\frac{1}{e}\right)^{\frac{K}{n} \ln 2}$$

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h bits one}] \stackrel{\text{Random}}{\approx} \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{\frac{hn}{K}} h}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[\text{h bits one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h = \left(1 - e^{\left(\frac{K}{n} \ln 2\right) \frac{n}{K}}\right)^{\frac{K}{n} \ln 2} = \left(\frac{1}{2}\right)^{\frac{K}{n} \ln 2}$$

with optimal h
(false positive)

Call this ϵ

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[\text{h bits one}] \stackrel{\text{Random}}{\approx} \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-e^{\frac{hn}{K}} h}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[\text{h bits one}] \stackrel{\text{Random}}{\approx} \left(1 - e^{\frac{hn}{K}}\right)^h = \left(1 - e^{\left(\frac{K}{n} \ln 2\right) \frac{n}{K}}\right)^{\frac{K}{n} \ln 2} = \left(\frac{1}{2}\right)^{\frac{K}{n} \ln 2}$$

with optimal h
(false positive)

Call this ϵ

$$\epsilon = 2^{-\frac{K}{n} \ln 2}$$

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[h \stackrel{\text{Random}}{\text{Bits}} \text{ one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-\frac{hn}{K}}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[h \stackrel{\text{Random}}{\text{Bits}} \text{ one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h = \left(1 - e^{\left(\frac{K}{n} \ln 2\right) \frac{K}{n}}\right)^{\frac{K}{n} \ln 2} = \left(\frac{1}{2}\right)^{\frac{K}{n} \ln 2}$$

with optimal h
(false positive)

Call this ϵ

$$\epsilon = 2^{-\frac{K}{n} \ln 2}$$
$$\ln \epsilon = (\ln 2) \left(-\frac{K}{n} \ln 2\right)$$

False Positive Rate

$$P[\text{single bit zero}] = \left(1 - \frac{1}{K}\right)^{hn} = \left(1 - \frac{1}{K}\right)^{K \cdot \frac{hn}{K}} \approx e^{\frac{hn}{K}}$$

$$P[h \stackrel{\text{Random}}{\text{Bits}} \text{ one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h \approx e^{-\frac{hn}{K}}$$

This is minimized when $h = \frac{K}{n} \ln 2$

$$P[h \stackrel{\text{Random}}{\text{Bits}} \text{ one}] \approx \left(1 - e^{\frac{hn}{K}}\right)^h = \left(1 - e^{\left(\frac{K}{n} \ln 2\right) \frac{n}{K}}\right)^{\frac{K}{n} \ln 2} = \left(\frac{1}{2}\right)^{\frac{K}{n} \ln 2}$$

with optimal h
(false positive)

Call this ϵ

$$\epsilon = 2^{-\frac{K}{n} \ln 2}$$

$$\ln \epsilon = (\ln 2) \left(-\frac{K}{n} \ln 2\right)$$

$$\frac{K}{n} = \frac{-1}{\ln 2} \ln \epsilon$$

Summary

Choose $\frac{K}{n} = -\frac{1}{\ln 2} \ln E$

↑
Bits per item

Choose $h = \frac{K}{n} \ln 2 = -\frac{1}{\ln 2} \ln E$

↑
number of hash functions

↓ error rate

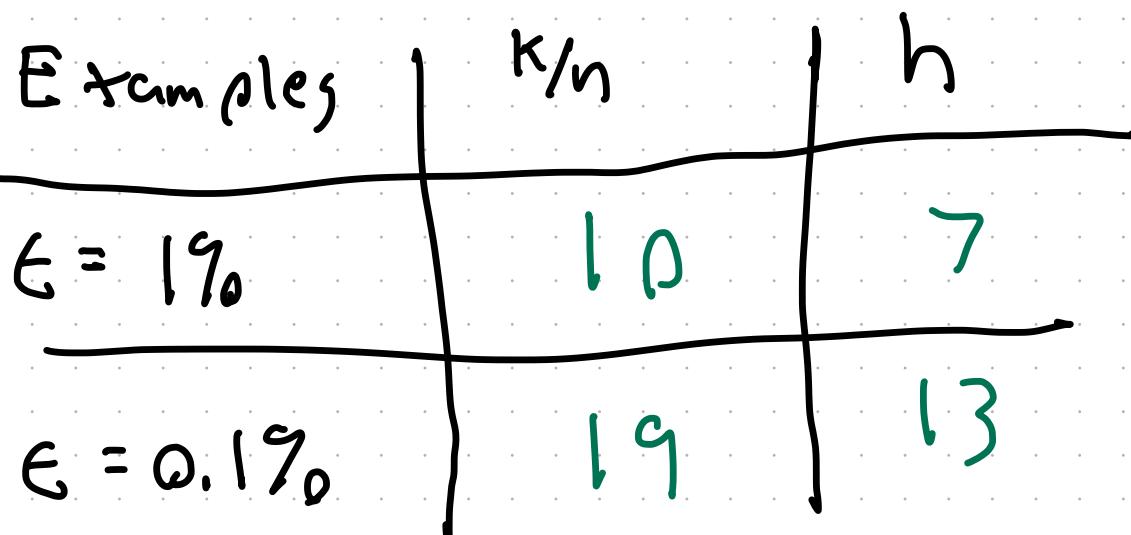
Summary

Choose $\frac{K}{n} = -\frac{1}{\ln 2} \ln \epsilon$

↑
bits per item

Choose $h = \frac{K}{n} \ln 2 = -\frac{1}{\ln 2} \ln \epsilon$

↑
number of hash functions



↓ error rate

URL Example

- 1 million URL's, 70 bytes on average
- 70 MB of data

URL Example

- 1 million URL's, 70 bytes on average
 - 70 MB of data
- Bloom Filter 1% error $\frac{k}{n} = 10$
 - 1.2 MB of data