# Data Structures

## Bloom Filters

Design and Analysis
of Algorithms I

# Bloom Filters: Supported Operations

Raison d'être : fast Inserts and Lookups.

Comparison to Hash Tables :

Pros : more space efficient.

Cons : ① can't store an associated object

② no deletions

③ small false positive probability ( i.e., might say x has been Inserted even though it hasn't been )

# Bloom Filters: Applications

**Original:** early spellcheckers.

**Canonical:** list of forbidden passwords

**Modern:** network routers.

— limited memory, need to be super-fast

# Bloom Filter: Under the Hood

**Ingredients:** ① array of $n$ bits $\left(\text{So } \frac{n}{|S|} = \text{\# of bits per object in data set } S\right)$

② $k$ hash functions $h_1, \dots, h_k$ (k = small constant)

**Insert(x):** for $i = 1, 2, \dots, k$

set $A[h_i(x)] = 1$ (whether or not bit already set to 1)

**Lookup(x):** return TRUE $\iff$ $A[h_i(x)] = 1$ for every $i = 1, 2, \dots, k$.

**Note:** no false negatives. (if $x$ was inserted, Lookup(x) guaranteed to succeed)

**But:** false positive if all $k$ $h_i(x)$'s already set to 1 by other insertions.

# Heuristic Analysis

**Intuition**: should be a trade-off between space and error (false positive) probability.

**Assume**: [not justified] all $h_i(x)$'s uniformly random and independent (across different $i$'s and $x$'s).

**Setup**: $n$ bits, insert data set $S$ into bloom filter.

**Note**: for each bit of $A$, the probability it's been set to 1 is (under above assumption):

Tim Roughgarden

Under the heuristic assumption, what is the probability that a given bit of the bloom filter (the first bit, say) has been set to 1 after the data set S has been inserted?

○ $(1 - 1/n)^{k|S|}$ — prob 1st bit = 0

○ $1 - (1 - 1/n)^{k|S|}$ — prob 1st bit = 1

○ $(1/n)^{|S|}$

○ $(1 - 1/n)^{|S|}$
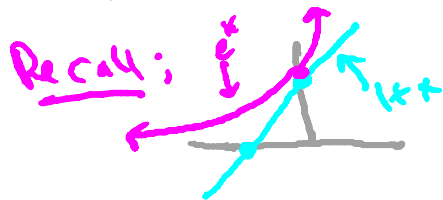
# Heuristic Analysis

**Intuition**: should be a trade-off between space and error (false positive) probability.

**Assume**: [not justified] all $h_i(x)$'s uniformly random and independent (across different $i$'s and $x$'s).

**Setup**: $n$ bits, insert data set $S$ into bloom filter.

**Note**: for each bit of $A$, the probability it's been set to $1$ is (under above assumption):

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

$$\leq 1 - e^{\frac{-k|S|}{n}} = 1 - e^{-k/b}$$

$\leftarrow b = \#$ of bits per object $(\frac{n}{|S|})$

**Recall**:



Tim Roughgarden

# Heuristic Analysis (con'd)

**Story So far:** probability a given bit is 1 is $\leq 1 - e^{-k/b}$

**So:** under assumption, for $x \notin S$, false positive probability is $\leq \boxed{[1 - e^{-k/b}]^k}$, where $b = \#$ of bits per object.

$\hookrightarrow$ error rate $\varepsilon$

**How to set k?:** For fixed $b$, $\varepsilon$ is minimized by setting $\boxed{k \approx (\ln 2) \cdot b}$

$\approx .693$

**Plugging back in:** $\varepsilon \approx \left(\frac{1}{2}\right)^{(\ln 2) b}$ (exponentially small in $b$)

or $b \approx 1.44 \log_2 \frac{1}{\varepsilon}$

**Ex:** with $b = 8$, choose $k = 5$ or $6$, error probability only $\approx 2\%$.

Tim Roughgarden