

The *Perfect Hash Filter*

A fast and efficient implementation of the *random approximate set* and *random approximate map* abstract data types.

Alexander Towell
atowell@siue.edu

Abstract

We define the semantics of the random approximate set and derive an immutable data structure that implements the *sstatic* random approximate set, denoted the Perfect Hash Filter. We analyze the theoretical properties of the Perfect Hash Filter, demonstrating that it is a *succinct* data structure that asymptotically obtains the information-theoretic lower bound on the average space complexity. Finally, we provide an implementation and compare it to the theoretical analysis.

Contents

1	Introduction	1
2	Approximate and oblivious sets	2
3	Perfect hash filters that trade time for space	3
4	A C++ implementation	3
4.1	Subtypes: dynamic polymorphism	3
5	Abstract data type of random approximate sets	5
6	Deterministic algorithm	5
6.1	Algebraic properties	6
6.2	Compatibility with the probabilistic model	6
7	Space complexity	8
7.1	Space efficiency of <i>unions</i> and <i>differences</i>	9

1 Introduction

The Perfect Hash Filter is a data structure that is suitable to implement two types of sets, *positive approximate sets*[?] and *oblivious sets*[?]. Informally, positive approximate sets generate false positives at specifiable rate and oblivious sets are a type of approximate set with additional *confidentiality* guarantees.

In ??, we precisely define concept of a set. Then, we define the approximate and oblivious positive sets. In ??, we derive the Perfect Hash Filter, a data structure that implements the approximate positive set and prove various properties, such as its absolute space efficiency. Following from that, in ?? we show how to tweak the implementation of the Perfect Hash Filter¹ to implement the *oblivious set*.

¹For instance, by varying the load factor.

2 Approximate and oblivious sets

A set is given by the following definition.

Definition 2.1. *A set is an unordered collection of distinct elements from a universe of elements.*

A countable set is a *finite set* or a *countably infinite set*. A *finite set* has a finite number of elements. For example,

$$S = \{1, 3, 5\}$$

is a finite set with three elements. A *countably infinite set* can be put in one-to-one correspondence with the set of natural numbers. The cardinality of a set S is a measure of the number of elements in the set, denoted by

$$|S|. \quad (1)$$

The cardinality of a *finite set* is a non-negative integer and counts the number of elements in the set, e.g.,

$$|\{1, 3, 5\}| = 3.$$

The abstract data type of the immutable *positive* approximate set[?] is given by the following definition.

Definition 2.2. *The abstract data type of the positive approximate set over a countably infinite universe \mathcal{U} has the following operations defined:*

$$\in: \mathcal{P}(\mathcal{U}) \times \mathcal{U} \mapsto \{\text{true}, \text{false}\}. \quad (2)$$

Let an element that is selected uniformly at random from the universe \mathcal{U} be denoted by X . A set $S^+(\varepsilon)$ is a positive approximate set of a set S with a false positive rate ε if the following conditions hold:

- (i) S is a subset of S^+ . This condition guarantees that no false negatives may occur.
- (ii) If X is not a member of S , it is a member of S^+ with a probability ε ,

$$P[X \in S^+ \mid X \notin S] = \varepsilon. \quad (3)$$

The optimal space complexity of *countably infinite* positive approximate sets is given by the following postulate.

Postulate 2.1. *The optimal space complexity of a data structure implementing the positive approximate set over countably infinite universes is independent of the type of elements and depends only the false positive rate ε as given by*

$$-\log_2 \varepsilon \text{ bits/element}. \quad (4)$$

The immutable abstract data type of the *oblivious set*[?] is given by the following definition.

Definition 2.3. *Assuming that the only information about a set of interested S is given by another set \check{S} , \check{S} is an oblivious positive set of S with a false positive rate ε if the following conditions hold:*

- (i) \check{S} is a positive approximate set of S with a false positive rate ε .
- (ii) The false positives are uniformly distributed over $\mathcal{U} \setminus S$.
- (iii) There is no efficient way to enumerate the elements in \check{S} .²

²That is, the true positives and false positives.

(iv) Any estimator of the cardinality of \mathcal{S} may only be able determine an approximate upper and lower bound, where the uncertainty may be traded for space-efficiency.

The absolute space efficiency of a data structure implementing an oblivious positive set consisting of m true positives with a false positive rate ε and an entropy β is given by

$$E(\varepsilon, m, \beta) = E \left[\frac{-(m + X) \log_2 \varepsilon}{\ell(\text{make_cipher_set}(m + X, \varepsilon))} \right], \quad (5)$$

where

$$X \sim \text{DU}(0, 2^\beta - 1). \quad (6)$$

The relative efficiency of the *optimal* positive oblivious set with entropy β to the *optimal* positive approximate set ($\beta = 0$) has an expectation given by

$$\text{RE}(\cdot, m, \beta) = 2^{-\beta} \sum_{k=0}^{2^\beta-1} \left(1 + \frac{k}{m} \right)^{-1}. \quad (7)$$

For a fixed β , as $m \rightarrow \infty$ the relative efficiency goes to 1.

3 Perfect hash filters that trade time for space

$\text{ph}_S: \mathcal{U} \mapsto [0, \dots, k]$
 $f: \mathbb{N} \mapsto \{0, 1, \dots, N\}^{[t]}$
 Compose two functions
 $\text{ph}_S: \mathcal{U} \mapsto \{0, 1, \dots, N\}^{[t]} = f \circ \text{ph}_S$

4 A C++ implementation

Random positive approximate sets are a well-defined concept. Any data structure T and set of functions dependent upon T required by the concept is an implementation of the model. Any generic algorithm (generic programming) parameterized by R which assumes R models a random approximate set may be applied to the data structure T .

We show a simple implementation.

4.1 Subtypes: dynamic polymorphism

TODO: break up each function/class/maybe method into separate listings and then talk about each separately.

Open set of data types, closed set of operations.

We may also use *dynamic polymorphism*. A dynamic polymorphic C++ interface of the *random approximate set* abstract data type is detailed next.

Assuming the elements of an approximate set are *non-enumerable*, i.e., only membership tests may be performed using the member-of $\in: \mathcal{U} \times \mathcal{P}(\mathcal{U})$ interface, set-theoretic operations may be implemented by storing each approximate set *as-is* and composing them together. For example, if we are given two approximate sets \mathcal{S}^1 and \mathcal{S}^2 , the union $\mathcal{S}^1 \cup \mathcal{S}^2$ is implemented by performing membership tests on both \mathcal{S}^1 and \mathcal{S}^2 and returning **true** if either membership test returns **true**, i.e.,

$$x \in (\mathcal{S}^1 \cup \mathcal{S}^2) \equiv (x \in \mathcal{S}^1) \vee (x \in \mathcal{S}^2). \quad (8)$$

Any set-theoretic composition may be implemented as a combination of unions and complements as described in ???. A C++ implementation for taking the union of two approximate sets is given by ??? and a C++ implementation for taking the complement of an approximate set is given by ???.

The following C++ implementation of a Perfect Hash Filter of type X models a first-order rate-distorted set R_ε^ω of type X where $\varepsilon = 2^{-K}$ and ω is the product of ε and the possible rate-distortion of the (rate-distorted) perfect hash function.

```
template <
    // Expected false positive rate  $\varepsilon = 2^{-K}$ .
    int K,
    // Models a function of the type  $\text{string\_view} \mapsto \text{int}$ .
    template PerfectHashFn,
    bool Complement = false
>
class perfect_hash_filter
{
public:
    template <typename Set>
    perfect_hash_filter(
        Set A,    // A is objective set to approximate
        float r, // r is load factor
        uint32_t seed)
    {
        h = make_hash_fn(seed);
        ph_ = build_perfect_hash_fn(h_, A, r);

        auto N = ceil(size(A) / r);
        hs_.resize(N);
        for (auto x : A)
            hs_[ph_(seed, x)] = h_(x) % K;
    }

    auto contains(X x) const
    {
        return hs_[ph_(x)] == h_(x) % K;
    }

private:
    PerfectHashFn<hash_fn> ph_;
    hash_fn h_;
    vector<uint32_t> hs_;
};

template <int K, typename PerfectHashFn>
auto fpr(perfect_hash_filter<K, PerfectHashFn, false>)
{
    return power<K, 2>{}; //  $2^{-K}$ 
}
```

```

template <int K, typename PerfectHashFn>
auto fnr(perfect_hash_filter<K, PerfectHashFn, false> A)
{
    return rate_distortion(A.hash_fn()*(1-fpr(A)));
}

template <int K, typename PerfectHashFn>
auto fpr(perfect_hash_filter<K, PerfectHashFn, true>)
{
    return rate_distortion(A.hash_fn()*(1-fpr(A))); ???
}

template <int K, typename PerfectHashFn>
auto fnr(perfect_hash_filter<K, PerfectHashFn, true>) { return power<-K,2>{}; }

```

The complement unary operator is given next.

```

template <typename K, bool Complement>
perfect_hash_filter<K, !Complement> operator~(perfect_hash_filter<K, Complement> a)
{
    return perfect_hash_filter<K, !Complement>(a);
}

```

The perfect hash filter, as a *first-order* rate-distorted set, is closed under complements, i.e., it remains a first-order approximation after complementation.

If Container models the concept of a PackedContainer, Container<Integer> requires $O(m/r \log 2N)$ bits.

5 Abstract data type of random approximate sets

A *type* is a set and the elements of the set are called the *values* of the type. An *abstract data type* is a type and a set of operations on values of the type. For example, the *integer* abstract data type is defined by the set of integers and standard operations like addition and subtraction.

A *data structure* is a particular way of organizing data and may implement one or more abstract data types. Two different data structures that implement the abstract data type should be *exchangable* in a computer program without changing the outward behavior of the program.

See section 4 to see how a C++ interface for the random approximate set data type may be defined. A popular implementation of an algorithm and data structure that generates approximate sets consistent with the model is known as the *Bloom filter*.

6 Deterministic algorithm

By ????, each independent observation of S^\pm generates an uncertain number of false positives and false negatives. In practice, this may or may not be the case.

Suppose we have a deterministic algorithm that generates approximate sets with a specified expected false positive rate ε and false negative rate ω . The algorithm, being deterministic, necessarily generates a certain set \mathcal{X} given an objective set \mathcal{S} . That is to say, the algorithm is a *total function*, $f: \mathcal{P}(\mathcal{U}) \mapsto \mathcal{P}(\mathcal{U})$, with an *image*

$$\text{image}(f) = \{ f(\mathcal{A}) \mid \mathcal{A} \in \mathcal{P}(\mathcal{U}) \} \subseteq \mathcal{P}(\mathcal{U}). \quad (9)$$

Since the function f may map more than one input to the same output, and some sets in the codomain may not be mapped to by any set in the domain, f is (possibly) a non-surjective, non-injective function.

6.1 Algebraic properties

The only thing we can say with certainty *a priori*³ about the image of f is that its members are subsets of \mathcal{U} and it contains the empty set \emptyset and universal set \mathcal{U} .

Definition 6.1. A σ -algebra is closed under...

The image of f is not necessarily a σ -algebra. However, the subsets of \mathcal{U} that may be constructed by countable complements, unions, and intersections for elements of the image is by definition a σ -algebra and is denoted by $\sigma(f)$.

Since $\sigma(f)$ is a set of sets closed under unions, intersections, and complements, it is a Boolean algebra defined by the six-tuple $(\sigma(f), \cup, \cap, \bar{\cdot}, \emptyset, \mathcal{U})$, e.g., set-theoretic operations over the above Boolean algebra are of the form

$$\sigma(f) \times \sigma(f) \mapsto \sigma(f). \quad (10)$$

Note that neither *positive* nor *negative* approximate sets are closed under complementation and thus are not Boolean algebras, but rather semi-rings under unions and intersections.

Suppose we have two distinct Boolean algebras, $(\sigma(f), \cup, \cap, \bar{\cdot}, \emptyset, \mathcal{U})$ and $(\sigma(g), \cup, \cap, \bar{\cdot}, \emptyset, \mathcal{U})$, where f and g are distinct deterministic algorithms that generate approximate sets for $\mathcal{P}(\mathcal{U})$. Set-theoretic operations over both Boolean algebras is the Boolean algebra $(\Sigma_{f \cup g}, \cup, \cap, \bar{\cdot}, \emptyset, \mathcal{U})$, where $\Sigma_{f \cup g} = \sigma(\sigma(f) \cup \sigma(g))$. Note that $\sigma(f), \sigma(g) \subseteq \Sigma_{f \cup g}$, so if we take the limit $\lim_{n \rightarrow \infty} \Sigma_{f_1 \cup \dots \cup f_n}$, where f_1, \dots, f_n are distinct mappings, $\Sigma_{f_1 \cup \dots \cup f_n}$ converges to $\mathcal{P}(\mathcal{U})$.

In practice, it is often trivial to implement a large family of random approximate set generators with distinct Boolean algebras given a single implementation, e.g., “randomly” seeding a Bloom filter’s hash function. As the size of this family goes to infinity, “randomly” selecting one of them each time we generate an approximate set generates a random sample consistent with the probability distribution given by ??.

6.2 Compatibility with the probabilistic model

How do we reconcile the fact that the algorithm is a *function* with the *probabilistic model*? In this context, the notion of *probability* quantifies our *ignorance*:

1. Given an objective set \mathcal{S} , we do not have complete *a priori* knowledge about the set \mathcal{X} it maps to (the approximate set model only provides *a priori* knowledge about the probability distribution \mathcal{S}^\pm). We acquire *a posteriori* knowledge⁴ by applying the algorithm and *observing* \mathcal{X} as its output.
2. Given a set \mathcal{X} that approximates some unknown objective set \mathcal{S} , we do not have complete *a priori* knowledge about \mathcal{S} . According to the probabilistic model, the only *a priori* knowledge we have is given by the specified *expected* false positive and false negative rates. We may acquire a posteriori knowledge by applying the algorithm to each set in $\mathcal{P}(\mathcal{U})$ and remembering the sets that map to \mathcal{X} .⁵ However, since f is (possibly) non-injective, one or more sets may map to \mathcal{X} and thus this process may not completely eliminate uncertainty. Additionally, the domain $\mathcal{P}(\mathcal{U})$ has a cardinality $2^{|\mathcal{U}|}$ and thus exhaustive maps are impractical to compute even for relatively small domains.⁶ However, we may still reduce our uncertainty by mapping some subset of interest.

³A priori knowledge is independent of experience.

⁴A posteriori knowledge is dependent on experience.

⁵If the approximate set is the result of the union, intersection, and complement of two or more approximate sets, then we must consider the closure.

⁶In the case of *countably infinite* domains, it is not even theoretically possible.

Suppose \mathcal{U} is finite. There are $2^{|\mathcal{U}|}$ subsets of \mathcal{U} , and each set in $\mathcal{P}(\mathcal{U}) \setminus \{\emptyset, \mathcal{U}\}$ may map (assuming both false positives and false negatives are possible) to *any* of subset. There are a total of

$$2^{2(|\mathcal{U}|-1)} \quad (11)$$

possible functions that map perform this mapping.⁷ Each of these functions is compatible with the probabilistic model. For instance, a Bloom filter (positive approximate set) may have a family of hash function that, for a particular binary coding of the elements of a given universal set, maps *every* element in the universal set to the same hash. Thus, for instance, no matter the objective set $\mathcal{X} \subseteq \mathcal{U}$, it will map to \mathcal{U} . The Bloom filter had a theoretically sound implementation, but only after empirical evidence was it discovered that it was not suitable. This is an extremely unlikely outcome in the case of large universal sets, but as the cardinality of the universal set decreases, the probability of such an outcome increases. Indeed, at $|\mathcal{U}| = 2$, the probability of this outcome is ?.

Thus, *a priori* knowledge, e.g., a theoretically sound algorithm, is not in practice sufficient (although for large universal sets, the probability is negligible). The suitability of an algorithm can only be determined by acquiring *a posteriori* knowledge.

We could explore the space of functions in the family and only choose those which, on some sample of objective sets of interest, generates the desired expectations for the false positive and false negative rates with the desired variances. Most of them will if constructed in the right sort of way.

A family of functions that are compatible with the probabilistic model is given by observing a particular realization $\mathcal{X} = \mathcal{S}^\pm$ and outputting \mathcal{X} on subsequent inputs of \mathcal{S} , i.e., caching the output of a *non-deterministic* process that conforms to the probabilistic model. This is essentially how well-known implementations like the Bloom filter work, where the pseudo-randomness comes from mechanical devices like hash functions that approximate random oracles.

The false positive rate of the approximate set corresponding to objective set \mathcal{X} is given by

$$\hat{\varepsilon}(\mathcal{X}) = \frac{1}{n} \sum_{x \in \bar{\mathcal{X}}} \mathbb{1}_{f(\mathcal{X})}(x), \quad (12)$$

where $n = |\bar{\mathcal{X}}|$.

Let \mathcal{U}_p denote the set of objective sets with cardinality p . The *mean* false positive rate,

$$\bar{\varepsilon} = \frac{1}{|\mathcal{U}_p|} \sum_{\mathcal{X} \in \mathcal{U}_p} \hat{\varepsilon}(\mathcal{X}), \quad (13)$$

is an unbiased estimator of ε and the population variance

$$s_\varepsilon^2 = \frac{1}{|\mathcal{U}_p|} \sum_{\mathcal{X} \in \mathcal{U}_p} \hat{\varepsilon}(\mathcal{X}), \quad (14)$$

is an unbiased estimator of $V[\mathcal{E}_n] = \varepsilon(1 - \varepsilon)/n$.

Proof. We imagine that the function f caches the output of a *non-deterministic* process that conforms to the probabilistic model. Thus, each time the function maps an objective set \mathcal{X} of cardinality p to its approximation, the algorithm *observes* a realization of $\mathcal{E}_n = \hat{\varepsilon}$. Thus,

$$\bar{\varepsilon} = \frac{1}{|\mathcal{U}_p|} \sum_{\mathcal{X}_j \in \mathcal{U}_p} \hat{\varepsilon}(\mathcal{X}_j) \quad (a)$$

$$= \frac{1}{|\mathcal{U}_p|} \sum_{\mathcal{X}_j \in \mathcal{U}_p} E[\mathcal{E}_n^{(j)}] = \varepsilon. \quad (b)$$

⁷There are a countably infinite number of algorithms that may carry out the computations that implement any particular function.

□

7 Space complexity

If the finite cardinality of a universe is u and the set is *dense* (and the approximation is also dense, i.e., the false negative rate is relatively small), then

$$O(u) \text{ bits} \quad (15)$$

are needed to code the set, which is independent of m , the false positive rate, and the false negative rate.

The lower-bound on the *expected* space complexity of a data structure implementing the *approximate set* abstract data type where the elements are over a *countably infinite* universe is given by the following postulate.

Postulate 7.1. *The information-theoretic lower-bound of a data structure that implements the countably infinite approximate set abstract data type has an expected bit length given by*

$$-(1 - \omega) \log_2 \varepsilon \text{ bits/element}, \quad (16)$$

where $\varepsilon > 0$ is the false positive rate and ω is the false negative rate.

The *relative space efficiency* of a data structure X to a data structure Y is some value greater than 0 and is given by the ratio of the bit length of Y to the bit length of X ,

$$\text{RE}(X, Y) = \frac{\ell(Y)}{\ell(X)}, \quad (17)$$

where ℓ is the bit length function. If $\text{RE}(X, Y) < 1$, X is less efficient than Y , if $\text{RE}(X, Y) > 1$, X is more efficient than Y , and if $\text{RE}(X, Y) = 1$, X and Y are equally efficient. The absolute space efficiency is given by the following definition.

Definition 7.1. *The absolute space efficiency of a data structure X , denoted by $\text{E}(X)$, is some value between 0 and 1 and is given by the ratio of the bit length of the theoretical lower-bound to the bit length of X ,*

$$\text{E}(X) = \frac{\ell}{\ell(X)}, \quad (18)$$

where $\ell(X)$ denotes the bit length of X and ℓ denotes the bit length of the information-theoretic lower-bound. The closer $\text{E}(X)$ is to 1, the more space-efficient the data structure. A data structure that obtains an efficiency of 1 is optimal.⁸

The *absolute* space efficiency of a data structure X implementing an approximate set \mathcal{S}^\pm with a false positive rate ε and false negative rate ω is given by

$$\text{E}(X) = \frac{-m(1 - \omega) \log_2 \varepsilon}{\ell(X)}, \quad (19)$$

where $m = |\mathcal{S}|$. The most useful sort of asymptotic optimality is with respect to the parameter m .

A well-known implementation of countably infinite *positive approximate set* is the Bloom filter[?] which has an expected space complexity given by

$$-\frac{1}{\ln 2} \log_2 \varepsilon \text{ bits/element}, \quad (20)$$

⁸Sometimes, a data structure may only obtain the information-theoretic lower-bound with respect to the limit of some parameter, in which case the data structure *asymptotically* obtains the lower-bound with respect to said parameter.

thus the absolute efficiency of the Bloom filter is $\ln 2 \approx 0.69$. A practical implementation of the *positive approximate set* abstract data type is given by the *Perfect Hash Filter*[?], which compares favorably to the Bloom filter in many circumstances.

The *Singular Hash Set*[?] is a *theoretical* data structure that optimally implements the abstract data types of the *approximate set* and the *oblivious set*[?] abstract data types.

7.1 Space efficiency of unions and differences

As a way to implement *insertions* and *deletions*, we consider the space efficiency of the set-theoretic operations of unions and differences of approximate sets.

Let $S_1 = \{x_{j_1}, \dots, x_{j_m}\}$ and suppose we wish to insert the elements x_{k_1}, \dots, x_{k_p} into S_1 . If X_1 is a mutable object, then an *insertion* operator may be applied on X_1 for each $x_{k_i}, i = 1, \dots, p$.

Alternatively, if X_1 is immutable, then we may construct another object, X_2 , that implements the set $S_2 = \{x_{k_1}, \dots, x_{k_p}\}$, and then apply the union function,

$$X_1 \cup X_2. \quad (21)$$

If we replace X_1 and X_2 by objects that implement *positive approximate sets* of S_1 and S_2 respectively, then by ??, the false positive rate of the resulting approximate set is $\hat{\epsilon}_1 + \hat{\epsilon}_2 - \hat{\epsilon}_1 \hat{\epsilon}_2$.

The space efficiency of this positive approximate set is given by the following theorem.

Theorem 7.1. *Given two countably infinite positive approximate sets S_1^+ and S_2^+ respectively with false positive rates $\hat{\epsilon}_1$ and $\hat{\epsilon}_2$, the approximate set $S_1^+ \cup S_2^+$, which has an induced false positive rate $\hat{\epsilon}_1 + \hat{\epsilon}_2 - \hat{\epsilon}_1 \hat{\epsilon}_2$, has an expected upper-bound on its absolute efficiency given by*

$$E(\hat{\epsilon}_1, \hat{\epsilon}_2 | \alpha_1, \alpha_2) = \frac{\log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 - \hat{\epsilon}_1 \hat{\epsilon}_2)}{\alpha_1 \log_2 \hat{\epsilon}_1 + \alpha_2 \log_2 \hat{\epsilon}_2}, \quad (22)$$

where

$$\begin{aligned} 0 < \alpha_1 &= \frac{|S_1|}{|S_1 \cup S_2|} \leq 1, \\ 0 < \alpha_2 &= \frac{|S_2|}{|S_1 \cup S_2|} \leq 1, \\ 1 &\leq \alpha_1 + \alpha_2. \end{aligned} \quad (23)$$

As $\hat{\epsilon}_j \rightarrow 1$ or $\hat{\epsilon}_j \rightarrow 0$ for $j = 1, 2$, or $(\hat{\epsilon}_1, \hat{\epsilon}_2) \rightarrow (1, 1)$, the absolute efficiency goes to 0.⁹

Proof. The proof comes in two parts. First, we prove eq. (22), and then we prove the bounds on α_1 and α_2 given by eq. (23).

Let X and Y denote optimally space-efficient data structures that respectively implement positive approximate sets S_1^+ and S_2^+ with false positive rates $\hat{\epsilon}_1$ and $\hat{\epsilon}_2$. By ??, their union has an induced false positive rate given by

$$\hat{\epsilon}_1 + \hat{\epsilon}_2 + \hat{\epsilon}_1 \hat{\epsilon}_2. \quad (a)$$

The information-theoretic lower-bound of the approximate set of $S_1 \cup S_2$ with the above false positive rate is given by

$$-|S_1 \cup S_2| \log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 + \hat{\epsilon}_1 \hat{\epsilon}_2) \text{ bits}. \quad (b)$$

⁹As $(\hat{\epsilon}_1, \hat{\epsilon}_2) \rightarrow (0, 0)$, the absolute efficiency depends on the path taken. In most cases, it goes to 0.

Since we assume we only have X and Y and it is not possible to enumerate the elements in either, we must implement their union by storing and separately querying both X and Y . Since X and Y are optimal, $\ell(X) = -|S_1| \log_2 \hat{\epsilon}_1$ and $\ell(Y) = -|S_2| \log_2 \hat{\epsilon}_2$. Making these substitutions yields an absolute efficiency ???

$$E = \frac{|S_1 \cup S_2| \log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 + \hat{\epsilon}_1 \hat{\epsilon}_2)}{|S_1| \log_2 \hat{\epsilon}_1 + |S_2| \log_2 \hat{\epsilon}_2}. \quad (c)$$

Letting

$$\alpha_1 = \frac{|S_1|}{|S_1 \cup S_2|} \text{ and } \alpha_2 = \frac{|S_2|}{|S_1 \cup S_2|}, \quad (d)$$

we may rewrite eq. (c) as

$$\frac{\log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 + \hat{\epsilon}_1 \hat{\epsilon}_2)}{\alpha_1 \log_2 \hat{\epsilon}_1 + \alpha_2 \log_2 \hat{\epsilon}_2}. \quad (22 \text{ revisited})$$

In the second part of the proof, we prove the bounds on α_1 and α_2 as given by eq. (23). Both α_1 and α_2 must be non-negative since each is the ratio of two positive numbers (cardinalities). If $|S_1| \ll |S_2|$, then $\alpha_1 \approx 0$. If $S_1 \supset S_2$, then $\alpha_1 = 1$. A similar argument holds for α_2 . Finally,

$$\alpha_1 + \alpha_2 = \frac{|S_1| + |S_2|}{|S_1 \cup S_2|} \quad (e)$$

has a minimum value by assuming that S_1 and S_2 are disjoint sets (i.e., their intersection is the empty set), in which case

$$\alpha_1 + \alpha_2 = \frac{|S_1| + |S_2|}{|S_1| + |S_2|} = 1. \quad (f)$$

□

See ?? for a contour plot of the expected lower-bound as a function of $\hat{\epsilon}_1$ and $\hat{\epsilon}_2$. As $\hat{\epsilon}_1 \rightarrow 0$ or $\hat{\epsilon}_2 \rightarrow 0$, the efficiency goes to 0.

The lower-bound on the efficiency of the union of approximate sets is given by the following corollary.

Corollary 7.1.1. *Given two positive, non-enumerable approximate sets with false positive rates $\hat{\epsilon}_1$ and $\hat{\epsilon}_2$, their union is an approximate set that has an efficiency that is expected to be greater than the lower bound given by*

$$\min E(\hat{\epsilon}_1, \hat{\epsilon}_2) = \frac{\log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 + \hat{\epsilon}_1 \hat{\epsilon}_2)}{\log_2 \hat{\epsilon}_1 \hat{\epsilon}_2}. \quad (24)$$

Corollary 7.1.2. *If $\epsilon_1 = \epsilon_2 = \epsilon$, then the absolute efficiency is given by*

$$E(\hat{\epsilon}|\alpha) = \left(1 + \frac{\log_2(2 - \hat{\epsilon})}{\log_2 \hat{\epsilon}}\right) \left(1 - \frac{\alpha}{2}\right), \quad (25)$$

where

$$0 \leq \alpha = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|} \leq 1, \quad (26)$$

which is a monotonically decreasing function with respect to ϵ and α with limits given by $\lim_{\hat{\epsilon} \rightarrow 0} E(\epsilon) = 1$ and $\lim_{\epsilon \rightarrow 1} E(\epsilon) = 0$.

See fig. 1 for a graphic illustration.

Figure 1: Expected lower-bound on efficiency of the union of two approximate sets with the same false positive rate ε , neither of which can be enumerated.

