

The *Oblivious Set* abstract data type

Alexander Towell
atowell@siue.edu

Abstract

We define the semantics of abstract data type for the *oblivious set*. We demonstrate a theoretical data structure, denoted the *Singular Hash Set*, that provides an optimal implementation of the oblivious set with respect to *entropy* and *space complexity*. Finally, we show how to use the *Bloom filter* and *Perfect Hash Filter* to implement oblivious sets and compare them to each other and the theoretically optimal *Singular Hash Set*.

Contents

1	Introduction	1
2	Oblivious sets	2
2.1	Oblivious object types	2
3	Probabilistic model	5
3.1	Positives and negatives	5
4	Entropy	6
4.1	Positives and negatives	6
5	Set estimators	7
6	The <i>Singular Hash Set</i>	9
6.1	Space complexity	12
6.2	Entropy	16
	Appendices	18
A	Probability mass of random bit length	18

1 Introduction

The *oblivious set*[?] is a fundamental data structure that may be used to construct other oblivious object types, like secure indices for Boolean Encrypted Search[?]

An *oblivious set* is an oblivious object type for the abstract data type of the *set*. The following operations are defined:

1. Approximate membership tests of specific elements may be performed with a false positive rate ε and a false negative rate ω . Note that there is no way to efficiently iterate over the elements.

2. The cardinality may be estimated to be within some range. The the degree of uncertainty can be made arbitrarily large entropy at the expense of its space complexity.
3. Set-theoretic operations like union, intersection, and complement generate oblivious sets that approximate the true operation as a function of the false positive and false negative rates and the degree of similarity between the exact sets under consideration.

In section 2, we precisely define the oblivious set.

In section 3, we derive the probabilistic model of *oblivious sets*. In section 4, we derive the *entropy* of *oblivious sets*. In section 5, we derive estimators of properties of *oblivious sets*. In section 6, we provide a theoretically optimal implementation of the oblivious set.

2 Oblivious sets

A set is given by the following definition.

Definition 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,

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

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

$$\mathbb{N} = \{1, 2, 3, 4, 5, \dots\}. \quad (1)$$

The cardinality of a set \mathcal{S} is a measure of the number of elements in the set, denoted by

$$|\mathcal{S}|. \quad (2)$$

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

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

Informally, an oblivious set of \mathcal{S} , denoted by $\check{\mathcal{S}}$, provides a *confidential* in-place binary representation such that very little information about \mathcal{S} is disclosed. To be a minimally *useful*, $\check{\mathcal{S}}$ must permit *membership tests* with respect to \mathcal{S} with specifiable false positive and false negative rates. In what follows, we provide a formal specification of the abstract data type of the oblivious set.

2.1 Oblivious object types

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.

Suppose we have an abstract data type denoted by T with a set of operators $\mathcal{F} = \{f_1, \dots, f_n\}$ that are (at least partially) functions of T . We denote that an *object* x in computer memory implements the abstract data type T by $T(x)$.

An *oblivious* object type[2] that implements the abstract data type T is a related type denoted by \check{T} that provides guarantees about what can be learned about an object $\check{T}(x)$ by looking at the binary representation of \check{x} .

Optimally, the only information that can be learned about \tilde{x} is given by the well-defined *behavior* of the the abstract data type on the set of operators \mathcal{F} . For instance, say an operator $g: [T] \mapsto \{\mathbf{true}, \mathbf{false}\}$ is defined but not in F , and $g(x) = \mathbf{true}$ for a particular object $T(x)$. Then, if the only information we have about x is given by $\tilde{x} = \tilde{T}(x)$, then $P[g(\tilde{x}) = g(x)] = 0.5$, i.e., we can do no better than a random guess.

The *oblivious set* is an abstract data type which *confidentially* approximates sets with two types of errors, false positives and false negatives. Thus, the oblivious part consists of two parts. First, it is a type of *approximate set*[1]. Second, it provides additional confidentiality guarantees.

The abstract data type of the immutable approximate set[1] is given by the following definition.

Definition 2. *The abstract data type of the approximate set over a universe \mathcal{U} has values given by the set $\mathcal{P}(\mathcal{U})$. At a minimum, a set must provide some way to test whether particular elements in \mathcal{U} are members of a particular set,*

$$\in: \mathcal{U} \times \mathcal{P}(\mathcal{U}) \mapsto \{\mathbf{true}, \mathbf{false}\}, \quad (3)$$

$$(4)$$

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

(i) *If X is a member of \mathcal{S} , it is not a member of \mathcal{S}^* with a probability ω ,*

$$P[X \notin \mathcal{S}^* \mid X \in \mathcal{S}] = \omega. \quad (5)$$

(ii) *If X is not a member of \mathcal{S} , it is a member of \mathcal{S}^+ with a probability ε ,*

$$P[X \in \mathcal{S}^* \mid X \notin \mathcal{S}] = \varepsilon. \quad (6)$$

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

Postulate 1. *The optimal space complexity of a data structure implementing the approximate set over a countably infinite universe is independent of the type of elements and depends only the false positive rate ε and false negative rate ω as given by*

$$-(1 - \omega) \log_2 \varepsilon \text{ bits/element}. \quad (7)$$

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

Definition 3 (Oblivious set). *Assuming that the only information about a set of interested $\mathcal{S} \subset \mathcal{U}$ is given by another set $\tilde{\mathcal{S}}$, $\tilde{\mathcal{S}}$ is an oblivious set of a set \mathcal{S} if the following conditions are hold:*

(i) *There is no efficient way to enumerate the elements in $\tilde{\mathcal{S}}^*$.¹*

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

Definition 4 (Approximate oblivious set). *Assuming the conditions specified for oblivious sets hold in definition 3, an oblivious set $\tilde{\mathcal{S}}^*$ is an approximate oblivious set of \mathcal{S} with a false positive rate ε and a false negative rate ω if the following additional conditions hold:*

¹That is, the true positives and false positives.

1. Each negative element tests positive with a probability ε and tests negative with a probability $1 - \varepsilon$. That is, each test is Bernoulli distributed, which is the maximum entropy distribution given that the false positive rate is ε .²

Assuming we do not have access to \mathcal{S} , the most accurate prediction possible when predicting whether an element is negative is ε .

$$P[X \in \mathcal{S}] = P[X \text{ is a false negative or } X \text{ is a true positive}]. \quad (8)$$

2. Each positive element tests negative with a probability ω and tests positive with a probability $1 - \omega$. That is, each test is Bernoulli distributed, which is the maximum entropy distribution given that the false negative rate is ω .³
3. By items 1 and 2, $\tilde{\mathcal{S}}^*$ is an approximate set[1] of \mathcal{S} with a false positive rate ε and false negative rate ω .

A *oblivious positive set* is a special case given by the following definition.

Definition 5. An oblivious set $\tilde{\mathcal{S}}^*$ with a false negative rate equal to zero is a *oblivious positive set* denoted by $\tilde{\mathcal{S}}^+$. By this definition, $\tilde{\mathcal{S}}^+$ is a superset of \mathcal{S} .

The *complement* of a *oblivious positive set* is given by the following definition.

Definition 6. An oblivious set $\tilde{\mathcal{S}}$ with a false positive rate equal to zero is a *oblivious negative set* denoted by $\tilde{\mathcal{S}}^+$. By this definition, $\tilde{\mathcal{S}}^+$ is a subset of \mathcal{S} .

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

$$E(\varepsilon, \omega, m, \beta) = E \left[\frac{-(1 - \omega)(m + X) \log_2 \varepsilon}{\text{BL}(Y)} \right], \quad (9)$$

where

$$X \sim \text{DU}(0, 2^\beta - 1) \quad (10)$$

and Y is a function of the random variable X .

The *Singular Hash Set* is an *optimal* implementation of the oblivious set abstract data type.

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

$$\text{RE}(\cdot, m, \beta) = 2^{-\beta} \sum_{k=0}^{2^\beta-1} \left(1 + \frac{k}{m} \right)^{-1}. \quad (11)$$

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

See ?? to see how a C++ interface for the approximate set abstract data type may be defined.

²If the universe is finite and there are n negatives, the number of false positives is binomially distributed with a mean εn .

³If there are p positives, the number of false positives is binomially distributed with a mean εn .

3 Probabilistic model

The uncertain number of false positives is given by the following theorem.

Theorem 1. *Given a set \mathcal{S} with m positives from a universe of u elements, the number of false positives in an approximate set \mathcal{S}^* with a false positive rate ε is a random variable denoted by FP_m with a distribution given by*

$$\text{FP}_m \sim \text{BIN}(u - m, \varepsilon). \quad (12)$$

The number of false negatives is given by the following theorem.

Theorem 2. *Given a set \mathcal{S} with m positives, the number of false negatives with respect to an approximate set \mathcal{S}^* with a false negative rate ω is a random variable denoted by FN_m with a distribution given by*

$$\text{FN}_m \sim \text{BIN}(m, \omega). \quad (13)$$

The *expected* cardinality is given by the following theorem.

Theorem 3 (Cardinality). *Given a set \mathcal{S} of cardinality m from a universe of u elements, an approximate set \mathcal{S}^* has an expected cardinality given by*

$$u\varepsilon + m(1 - \varepsilon - \omega), \quad (14)$$

where ε is the false positive rate and ω is the false negative rate.

3.1 Positives and negatives

The distribution of false positives and false negatives are Bernoulli distributed random variables conditioned on a particular number of positives. The distribution of positives (and negatives) is given by the following definition.

Definition 7. *The number of positives in a universe of u elements is uncertain. We model the uncertainty as a discrete random variable, denoted by P , with a probability mass function⁴*

$$\text{pP}(p | u) \quad (15)$$

and a support $\{0, \dots, u\}$. Conversely, the distribution of negatives is a random variable $N = u - P$ with a probability mass function

$$\text{pN}(n | u) = \text{pP}(u - n | u). \quad (16)$$

The form the probability mass function $\text{pP}(\cdot)$ takes cannot be a priori specified, although it may be estimated with an empirical probability.

Modeling the distribution of positives provides a complete specification for the distribution of false positives and false negatives (and true positives and true negatives).

Example 1 The *expected* number of false positives is give by the expectation

$$\text{E}[\text{FP}] = \sum_{p=0}^u \sum_{f_p=0}^{u-p} f_p \cdot \text{pP}(p | u) \text{pFP}(f_p | p, u, \varepsilon) \quad (a)$$

$$= \sum_{p=0}^u \text{pP}(p | u) \text{E}[\text{FP}_p | u] = \sum_{p=0}^u \text{pP}(p | u) (u - p) \varepsilon \quad (b)$$

$$= \varepsilon \left(u - \sum_{p=0}^u m \text{pP}(p | u) \right) = \varepsilon (u - \text{E}[P]) . \quad (c)$$

⁴The probability mass function of a random variable X is denoted by $\text{pX}(\cdot)$.

Note that $E[N] = u - E[P]$, thus

$$E[FP] = \varepsilon E[N]. \quad (d)$$

The joint probability mass function of positives, false positives, and false negatives is given by

$$p(p, f_p, f_n \mid u, \varepsilon, \omega) = p_P(p \mid u) p_{FP}(f_p \mid p, u, \varepsilon) p_{FN}(f_n \mid p, u, \omega). \quad (17)$$

4 Entropy

Theorem 4. *The entropy of the uncertain number of false positives and false negatives is given by*

$$c + \frac{1}{2} \log_2((u - m)m\varepsilon(1 - \varepsilon)\omega(1 - \omega)) + \mathcal{O}\left(\frac{u}{m(u - m)}\right), \quad (18)$$

where $c = \log_2(2\pi e)$.

Proof. The entropy of the joint distribution of false positives and false negatives given that m are positive is given by

$$\mathcal{H}(FP_m, FN_m). \quad (a)$$

By ??, FP_m and FN_m are independent. Thus,

$$\mathcal{H}(FP_m, FN_m) = \mathcal{H}(FP_m) + \mathcal{H}(FN_m). \quad (b)$$

The entropy of FP_m is defined as

$$\mathcal{H}(FP_m) = - \sum_{f_p=0}^{u-m} \log_2 p_{FP_m}(f_p \mid u, \varepsilon) p_{FP_m}(f_p \mid u, \varepsilon). \quad (c)$$

$$\mathcal{H}(FP_m) = c + \log_2 \sqrt{u - m} + \log_2 \sqrt{\varepsilon} + \log_2 \sqrt{1 - \varepsilon} + \mathcal{O}\left(\frac{1}{u - m}\right), \quad (d)$$

$$\mathcal{H}(FN_m) = c + \log_2 \sqrt{m} + \log_2 \sqrt{\omega} + \log_2 \sqrt{1 - \omega} + \mathcal{O}\left(\frac{1}{m}\right), \quad (e)$$

and $c = \log_2 \sqrt{2\pi e}$. Summing these and simplifying yields the result

$$\mathcal{H}(FP_m, FN_m) = c + \frac{1}{2} \log_2((u - m)m\varepsilon(1 - \varepsilon)\omega(1 - \omega)) + \mathcal{O}\left(\frac{u}{m(u - m)}\right), \quad (f)$$

where $c = \log_2(2\pi e)$. □

4.1 Positives and negatives

The distribution of false positives and false negatives are Bernoulli distributed random variables conditioned on a particular number of positives. The distribution of positives (and negatives) is given by the following definition.

Definition 8. *The number of positives in a universe of u elements is uncertain. We model the uncertainty as a discrete random variable, denoted by P , with a probability mass function⁵*

$$p_P(p | u) \quad (19)$$

and a support $\{0, \dots, u\}$. Conversely, the distribution of negatives is a random variable $N = u - P$ with a probability mass function

$$p_N(n | u) = p_P(u - n | u). \quad (20)$$

The form the probability mass function $p_P(\cdot)$ takes cannot be a priori specified, although it may be estimated with an empirical probability.

Modeling the distribution of positives provides a complete specification for the distribution of false positives and false negatives (and true positives and true negatives).

Example 2 The *expected* number of false positives is give by the expectation

$$E[FP] = \sum_{p=0}^u \sum_{f_p=0}^{u-p} f_p \cdot p_P(p | u) p_{FP}(f_p | p, u, \varepsilon) \quad (a)$$

$$= \sum_{p=0}^u p_P(p | u) E[FP_p | u] = \sum_{p=0}^u p_P(p | u) (u - p) \varepsilon \quad (b)$$

$$= \varepsilon \left(u - \sum_{p=0}^u m p_P(p | u) \right) = \varepsilon (u - E[P]) . \quad (c)$$

Note that $E[N] = u - E[P]$, thus

$$E[FP] = \varepsilon E[N] . \quad (d)$$

The joint probability mass function of positives, false positives, and false negatives is given by

$$p(p, f_p, f_n | u, \varepsilon, \omega) = p_P(p | u) p_{FP}(f_p | p, u, \varepsilon) p_{FN}(f_n | p, u, \omega) . \quad (21)$$

Since FP and FN are independent, the joint entropy of P , FP, and FN is given by

$$\mathcal{H}(P, FP, FN) = \mathcal{H}(P) + \mathcal{H}(FP | P) + \mathcal{H}(FN | u - P) \quad (22)$$

$$= \mathcal{H}(P) + \sum_{p=0}^u \mathcal{H}(FP | p) + \sum_{n=0}^u \mathcal{H}(FN | n) . \quad (23)$$

5 Set estimators

Given an approximate set \mathcal{S}^* with a false positive rate ε and a false negative rate ω , the *method of moments* estimator of the cardinality of \mathcal{S} is given by

$$\hat{m} = \frac{\text{cardinality}(\mathcal{S}^*) - \varepsilon u}{1 - \varepsilon - \omega} . \quad (24)$$

⁵The probability mass function of a random variable X is denoted by $p_X(\cdot)$.

Since the optimal space complexity is $-\log_2 \varepsilon$ per element, any data structure that implements an approximate set with a false positive rate ε obtains the maximum entropy, i.e., the maximum entropy per element is given by

$$-\log_2 \varepsilon \text{ bits/element.} \quad (25)$$

Thus, the *entropy* of an implementation of the approximate set is given by

$$\frac{H(\mathcal{S}^+)}{m} = -\log_2 \varepsilon \quad (26)$$

Theorem 5. *An unbiased estimator of the cardinality of a countably infinite approximate set \mathcal{S}^+ of a set \mathcal{S} with a false positive rate ε is given by*

$$\hat{m} = \frac{\text{BL}(\mathcal{S}^+)}{b}, \quad (27)$$

where BL is the bit length function and b is the expected bits per element.

Proof. The expected bit length is given by

$$-mb$$

where m is the cardinality of \mathcal{S} . Thus, the *method of moments* estimator is given by assuming the bit length realizes the expected value,

$$\text{BL}(\mathcal{S}^+) = mb. \quad (a)$$

Solving for m results in the estimator

$$\hat{m} = \frac{\text{BL}(\mathcal{S}^+)}{b}. \quad (b)$$

□

Given an approximate set \mathcal{S}^* with a false positive rate ε and a false negative rate ω , the *method of moments* estimator of the cardinality of \mathcal{S} is given by

$$\hat{m} = \frac{|\mathcal{S}^*| - \varepsilon u}{1 - \varepsilon - \omega}. \quad (28)$$

NOTE: The exact oblivious set bit length reveals nothing about the size of set \mathcal{S} since it only depends on $|\mathcal{U}|$.

By Kerckhoffs's principle, we assume the algorithms are known. For instance, we assume the *space complexity* with respect to the cardinality of the *exact* set is known.

Thus, the *cardinality*, a unary function of $\check{\mathcal{S}}^*$, may be estimated by using the information about the *expected* bit length. If the expected bit length as a function of m is given by $f(m)$, where m is the cardinality of the exact set, then a *method of moments* estimator is given by assuming the bit length realizes the expected value,

$$\text{BL}(\check{\mathcal{S}}^*) = f(m), \quad (29)$$

and solving for m , resulting in the estimator

$$\hat{m} = f^{-1}(\text{BL}(\check{\mathcal{S}}^*)). \quad (30)$$

Consider the following example of the *Perfect Hash Filter*[?].

Example 3 The *optimal Perfect Hash Filter* has a space complexity given by

$$\text{BL}(\check{\mathcal{S}}^*) = m \log_2 \frac{e}{\varepsilon}. \quad (\text{a})$$

Solving for m results in the estimator

$$\hat{m} = \frac{\text{BL}(\check{\mathcal{S}}^*)}{\log_2 \frac{e}{\varepsilon}}. \quad (\text{b})$$

Depending on the entropy of the random bit length of the oblivious set object type, cardinality estimators with very low variance may be obtainable. Thus, in order to increase the uncertainty, we must artificially inflate the bit length of $\check{\mathcal{S}}^*$. One simple way of achieving this is to randomly sample a positive integer from 0 to N , and insert

6 The Singular Hash Set

In what follows, we provide a theoretical implementation of the *oblivious set* that obtains optimality in the following ways:

1. The space complexity obtains the theoretical lower-bound of an approximate set with a false positive rate ε and a false negative rate ω .
2. The entropy obtains the upper bound. This is necessarily the case since it obtains the optimal space complexity.

TODO: make an absolute efficiency measure for a given entropy level. Mainly, this means the entropy of the cardinality. You can trade entropy on this for space efficiency.

Definition 9 (Cartesian product). *Let $\mathcal{X}_1, \dots, \mathcal{X}_n$ denote n sets. The set $\mathcal{X}_1 \times \dots \times \mathcal{X}_n = \{(x_1, \dots, x_n) : x_1 \in \mathcal{X}_1 \wedge \dots \wedge x_n \in \mathcal{X}_n\}$ is called the Cartesian product of sets $\mathcal{X}_1, \dots, \mathcal{X}_n$.*

A shorthand notation for the Cartesian product $\mathcal{X} \times \mathcal{X} \times \mathcal{X}$ is denoted by \mathcal{X}^3 .

The binary set $\{0, 1\}$ is denoted by \mathcal{B} . The set of all bit (binary) strings of length n is therefore \mathcal{B}^n . The cardinality of \mathcal{B}^n is

$$|\mathcal{B}^n| = 2^n. \quad (31)$$

The set of all bit strings of length n or less is denoted by $\mathcal{B}^{\leq n}$, which has a cardinality $2^{n+1} - 1$. The countably infinite set of all bit strings, $\mathcal{B}^{\leq \infty}$, is also denoted by \mathcal{B}^* .

The bit length of an object x is denoted by

$$\text{BL}(x), \quad (32)$$

e.g., the bit length of any $x \in \mathcal{B}_n$ is $\text{BL}(x) = n$.

A (convenient) one-to-one correspondence between \mathcal{B} and \mathbb{N} is given by the following definition.

Definition 10. *Let the set of bit strings \mathcal{B} and the set of natural numbers \mathbb{N} have the bijection given by*

$$(b_1 b_2 \dots b_m) \longleftrightarrow 2^m + \sum_{j=1}^m 2^{m-j} b_j. \quad (33)$$

We denote the mapping of a bit string (or natural number) x by x' .

An important observation of this mapping is that a natural number n maps to a bit string n' of length $\text{BL}(n') = \lfloor \log_2 n \rfloor$.

Example 4 Consider the number 100 (base 10). To determine the bit string $100'$, we perform the following steps:

1. The length of $100'$ is $\lfloor \log_2 100 \rfloor = 6$.
2. Subtract $2^6 = 64$ from 100 which results in 36.
3. 36_{10} in binary is 100100_2 .
4. $100' = (100100)$.

A *hash function* is related to countable sets \mathcal{B} and \mathcal{B}_n and is given by the following definition.

Definition 11. A hash function $\mathbf{h}: \mathcal{B} \mapsto \mathcal{B}^n$ is a function such that all bit strings of arbitrary-length are mapped (hashed) to bit strings of fixed-length n . For a given $x \in \mathcal{B}$, $y = \mathbf{h}(x)$ is denoted the hash of x .

The Singular Hash Set is a data type that implements the *oblivious set* abstract data type. The implementation consists of a product data structure (tuple) $\mathcal{B}^k \times \mathcal{B}^*$, an algorithm that *generates* the data structure, and an algorithm that implements the *member-of* function by *querying* the data structure.

The Singular Hash Set assumes that random oracles are available.

Definition 12. A random oracle, denoted by $\mathbf{h}^*: \mathcal{B}^\infty \mapsto \mathcal{B}^\infty$, is a theoretical hash function whose output is uniformly distributed over the elements of \mathcal{B}^∞ .

The implementation of the algorithm that generates the data structure for the Singular Hash Set that implements an *approximate oblivious set* is given by the following theorem.

Theorem 6 (Singular Hash Set). *The data structure generated by*

$$\leftarrow^* \text{make_singular_hash_set}(\mathcal{S}, \varepsilon) \quad (34)$$

as given by algorithm 3 is an oblivious set of $\mathcal{S} \subset \mathcal{U}$ with a false positive rate $\varepsilon \in \{2^{-k} : k \in \mathbb{N}\}$ and a false negative rate $\omega \in \{\frac{j}{m} : j \in \{1, \dots, m\}\}$.

Proof. In order for the Singular Hash Set to be an *oblivious set* of $\check{\mathcal{S}}^*$, it must also be an *approximate set*. To be an approximate set of \mathcal{S} , it must satisfy the two conditions given by definition 2:

1. \mathcal{S} is a subset of $\check{\mathcal{S}}^*$. This condition guarantees that no *false negatives* may occur.
2. An element in \mathcal{B} that is not a member of \mathcal{S} is a member of $\check{\mathcal{S}}^*$ with a probability ε , denoted the false positive rate, i.e.,

$$\mathbb{P}[\mathbf{X} \in \check{\mathcal{S}}^* \mid \mathbf{X} \notin \mathcal{S}] = \varepsilon. \quad (\text{a})$$

To prove the first condition, note that algorithm 4 tests any element x for membership in \mathcal{S}^+ by computing the hash of x concatenated with the bit string b_n and returning **true** if the hash is h_k where algorithm 3 finds bit strings b_n and h_k such that each element of \mathcal{S} concatenated with b_n hashes to h_k .

To prove the second condition, suppose we have a set $\mathcal{S} = \{x_1, \dots, x_m\}$ and each element in \mathcal{S} hashes to $y = \mathbf{h}(x_1)$. By ??, $\mathbf{h}: \mathcal{B}^* \mapsto \mathcal{B}^k$ approximates a random oracle and thus uniformly distributes over its domain of 2^k possibilities. Since y is a particular element in \mathcal{B}^k , the probability that an element not in \mathcal{S} hashes to y is 2^{-k} . \square

Algorithm 1: Implementation of `make_singular_hash_set` over a finite universe \mathcal{U}

in : \mathcal{S} is a subset of a finite universal set \mathcal{U} .
out : An *oblivious* exact set of \mathcal{S} .

```
1 function make_singular_hash_set( $\mathcal{S}, \varepsilon, \omega$ )
2    $\mathcal{S}_B \leftarrow \{\text{encode}_{\mathcal{U} \rightarrow \mathcal{B}}(x) : x \in \mathcal{S}\}$ 
3    $\overline{\mathcal{S}}_B \leftarrow \{\text{encode}_{\mathcal{U} \rightarrow \mathcal{B}}(x) : x \in \overline{\mathcal{S}}\}$ 
   // To find the smallest bit string, search for a bit string of length  $n$  in ascending
   order.
4   for  $n \leftarrow 0$  to  $\infty$  do
5     for  $j \leftarrow 1$  to  $2^n$  do
6       found  $\leftarrow$  true
       // To maximize entropy we try bit strings of length  $n$  in random order.
7        $b_n \leftarrow$  a bit string of length  $n$  randomly drawn from  $\mathcal{B}_n$  without replacement
8        $h_k \leftarrow$  null
9       for  $x \in \mathcal{S}_B$  do
10        if  $h_k = \text{null}$  then
11           $h_k \leftarrow h(x \# b_n) \bmod (k+1)$ 
12        end
13        else if  $h \neq h(x \# b_n) \bmod (k+1)$  then
14          found  $\leftarrow$  false
15        end
16      end
17      for  $x \in \overline{\mathcal{S}}_B$  do
18        if  $h = h(x \# b_n) \bmod (k+1)$  then
19          found  $\leftarrow$  false
20        end
21      end
22      if found then
23        // This tuple is the data structure of the Singular Hash Set.
24        return  $(h_k, b_n)$ 
25      end
26    end
  end
```

Condition for maxium entropy:

The countably infinite intersection of separate instances of an approximate oblivious set of \mathcal{S} where $\omega > 0$ is the empty set,

$$\bigcap_j \check{\mathcal{S}}_j^* = \emptyset, \quad (35)$$

where $\check{\mathcal{S}}_1^*, \check{\mathcal{S}}_2^*, \dots$ are random instances of an approximate oblivious set of \mathcal{S} .

TODO: this directly follows from the sampling distribution and taking the limit.

In the case of a positive approximate oblivious set, where $\omega = 0$, the countably infinite intersection of separate instances of a positive approximate oblivious set of \mathcal{S} is \mathcal{S} ,

$$\bigcap_j \check{\mathcal{S}}_j^+ = \mathcal{S}. \quad (36)$$

A countable union of separate instances of a negative approximate oblivious set of \mathcal{S} is the complement of \mathcal{S}

$$\bigcap_j \check{\mathcal{S}}_j^+ = \overline{\mathcal{S}}. \quad (37)$$

Exact oblivious set:

$$k(\rho u - 1) + (u \log_2(1 - 2^{-k})(\rho - 1)) \quad (38)$$

$$k(\rho u - 1)/m + (u \log_2(1 - 2^{-k})(\rho - 1))/m \quad (39)$$

$$k \left(1 - \frac{1}{u\rho} \right) + \log_2 \left(1 - 2^{-k} \right) \left(1 - \frac{1}{\rho} \right) \quad (40)$$

$$k(\rho - \frac{1}{u}) + \log_2(1 - 2^{-k}(\rho - 1)) \text{ bits/element in universal set} \quad (41)$$

Theorem 7.

$$\min_k k(\rho - \frac{1}{u}) + \log_2(1 - 2^{-k}(\rho - 1)) \text{ bits/element in universal set} \quad (42)$$

The above algortihm for generating exact oblivious sets has a space complexity of u bits. Assuming m is not too much smaller than u , i.e., \mathcal{S} is dense, this is the theoretical lower-bound for sets over finite universes.

If m is not dense, then trying larger singular hash lengths will result in a better lower-bound.

Suppose you are interested in creating exact oblivious sets over the universe of bit strings up to length n , denoted by $\mathcal{B}_{\leq n}$. Then, there are $|\mathcal{B}_{\leq n}| = 2^{n+1} - 1$ bit strings in the universe. In this case, the set has a lower-bound given by $\mathcal{O}(n)$.

6.1 Space complexity

The probability that every element of \mathcal{S} collides for a particular bit string in ?? is given by the following theorem.

Theorem 8. *The probability that a bit string $b \in \mathcal{B}$ results in perfect collision in ?? is given by*

$$\varepsilon^{m-1}, \quad (43)$$

where $m = |\mathcal{S}|$ and ε is the specified false positive rate.

Algorithm 2: Implementation of `make_singular_hash_set`

in : \mathcal{S} is finite subset of a universal set \mathcal{U} , ε is the false positive rate, and ω is the maximum false negative rate.

out: An *oblivious, approximate set* \mathcal{S}^* .

```
1 function make_singular_hash_set( $\mathcal{S}, \varepsilon, \omega$ )
2    $\mathcal{S}_{\mathcal{B}} \leftarrow \{\text{encode}_{\mathcal{U} \rightarrow \mathcal{B}}(x) : x \in \mathcal{S}\}$ 
   //  $p$  is the minimum number of elements in  $\mathcal{S}$  that must be true positives.
3    $p \leftarrow (1 - \omega)|\mathcal{S}|$ 
   // A random oracle that hashes to  $k$  bits has a probability  $\varepsilon = 2^{-k}$  of hashing to a
   // particular value.
4    $k \leftarrow -\log_2 \varepsilon$ 
   // To find the smallest bit string, search for a bit string of length  $n$  in ascending
   // order.
5   for  $n \leftarrow 0$  to  $\infty$  do
6     for  $j \leftarrow 1$  to  $2^n$  do
7       found  $\leftarrow$  true
       // To maximize entropy we try bit strings of length  $n$  in random order.
8        $b_n \leftarrow$  a bit string of length  $n$  randomly drawn from  $\mathcal{B}_n$  without replacement
       // A minimum of  $p$  elements must have the same hash.
9        $\mathcal{M} \leftarrow$  all subsets of  $p$  elements from  $\mathcal{S}$ 
10      for  $\mathcal{P} \in \mathcal{M}$  do
        // Let the  $p$  elements in  $\mathcal{P}$  be denoted by  $x^{(1)}, \dots, x^{(p)}$ .
11         $h_k \leftarrow \mathbf{h}(x^{(1)} \uplus b_n) \bmod k$ 
12        for  $i \leftarrow 2$  to  $p$  do
13          if  $h_k \neq \mathbf{h}(x^{(i)} \uplus b_n) \bmod k$  then
14            found  $\leftarrow$  false
15          end
16        end
17      end
18      if found then
        // This tuple is sufficient to code the Singular Hash Set.
19        return  $(h_k, b_n)$ 
20      end
21    end
22  end
```

Algorithm 3: Implementation of `make_singular_hash_set`

in : \mathcal{S} is a *finite* set over \mathcal{U} , $\varepsilon \in \{2^{-k} : k \in \mathbb{N} \cup \{0\}\}$ is the false positive rate, and $\omega \in \{j/|\mathcal{S}| : j = 1, \dots, |\mathcal{S}|\}$ is the false negative rate.

out : An approximate *oblivious set* $\tilde{\mathcal{S}}^*$.

```
1 function make_singular_hash_set( $\mathcal{S}, \varepsilon, \omega$ )
2    $\mathcal{S}_{\mathcal{B}} \leftarrow \{\text{encode}_{\mathcal{U} \rightarrow \mathcal{B}}(x) : x \in \mathcal{S}\}$ 
   //  $p$  is the number of elements in  $\mathcal{S}$  that must be false negatives.
3    $p \leftarrow (1 - \omega)|\mathcal{S}|$ 
   // A random oracle that hashes to  $k$  bits has a probability  $\varepsilon = 2^{-k}$  of hashing to a
   // particular value.
4    $k \leftarrow -\log_2 \varepsilon$ 
   // To find the smallest bit string, search for a bit string of length  $n$  in ascending
   // order of length.
5   for  $n \leftarrow 0$  to  $\infty$  do
6     for  $j \leftarrow 1$  to  $2^n$  do
7       // To maximize entropy we try bit strings of length  $n$  in random order.
        $b_n \leftarrow$  a bit string randomly drawn from  $\mathcal{B}_n$  without replacement
8        $\text{def} = \frac{h(x)}{\mathbf{h}(x \uplus b_n) \bmod k}$  for  $x \in \mathcal{S}$ 
9       if  $|\mathcal{X}| = p$  then
10        // This tuple is sufficient to code the Singular Hash Set.
11        return  $(h_k, b_n)$ 
12      end
13    end
  end
```

Algorithm 4: Implementation of `contains`

in : \mathcal{S}^+ is the Singular Hash Set to query and x is the element to test for membership.

out : **true** if $x \in \mathcal{S}^+$ otherwise **false**. If we rephrase this with respect to $\mathcal{S} \subset \mathcal{S}^+$, then **true** if $x \in \mathcal{S}$ and otherwise **true** with probability ε and **false** with probability $1 - \varepsilon$, where ε is the false positive rate of \mathcal{S}^+ .

```
1 function contains( $\mathcal{S}^+, x$ )
   // The Singular Hash Set  $\mathcal{S}^+$  is coded by the tuple  $(b_n, h_k)$ , where  $b_n$  is a bit string of
   // length  $n$  and  $h_k$  is the singular hash.
2    $k = \text{BL}(h_k)$ 
3   if  $\mathbf{h}(x \uplus b_n) \bmod k = h_k$  then
4     // False positives occur with probability  $\varepsilon = 2^{-k}$ .
5     return true
6   else
7     return false
```

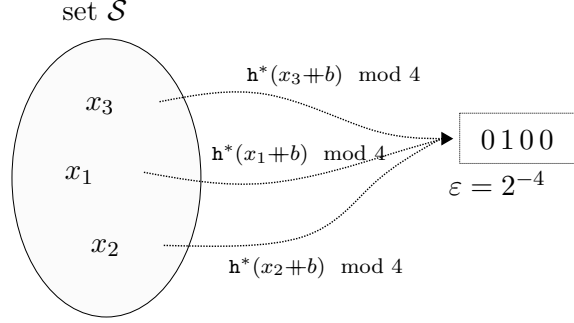


Figure 1: *Singular Hash Set* over a *countably infinite* universe

Proof. Suppose we have a set $\mathcal{S} = \{x_1, \dots, x_m\}$ and x_1 hashes to $y = \mathbf{h}_{\mathcal{S}}(x_1)$, where $\mathbf{h}_{\mathcal{S}}: \mathcal{B}^* \mapsto \mathcal{B}^k$ is a random hash function that uniformly distributes over its domain of 2^k possibilities. Since y is a particular element in \mathcal{B}^k , the probability that x_j for $j = 2, \dots, m$ hashes to y is given by

$$\frac{1}{2^k} = \varepsilon. \quad (\text{a})$$

Since $\mathbf{h}_{\mathcal{S}}$ is a random hash function, the hashes of x_1, \dots, x_m are independent. Thus, the joint probability that x_2, \dots, x_m hash to y is given by the product of their marginal probabilities

$$\varepsilon^{m-1}. \quad (\text{b})$$

□

Theorem 9. *The expected bit length of the Singular Hash Set obtains the information-theoretic lower-bound given by*

$$-\log_2 \varepsilon \text{ bits/element}, \quad (44)$$

where ε is the false positive rate.

Proof. The space required for the Singular Hash Set found by algorithm 3 is of the order of the length n of the bit string b_n . Therefore, for space efficiency, the algorithm exhaustively searches for a bit string in the order of increasing size n .

We are interested in the first case when a perfect collision occurs, which is a geometric distribution with probability of success ε^{m-1} as given by the discrete random variable

$$Q \sim \text{GEO}(\varepsilon^{m-1}). \quad (\text{a})$$

By definition 10, the n^{th} trial uniquely maps to a bit string of length $m = \lfloor \log_2 n \rfloor$. Thus, the bit string is a random length given approximately by

$$N = \log_2 Q \text{ bits}. \quad (\text{b})$$

This is a slight *overestimate* since we are simplifying by avoiding the floor function.

We approximate the logarithm with a second-order Taylor series around the *expected* value of Q as given by

$$N \approx \log_2 E[Q] - \frac{\log_2 e}{E[Q]} (Q - E[Q])^2 \text{ bits}. \quad (\text{c})$$

We are interested in the *expected* value of N ,

$$E[N] \approx \log_2 E[Q] - \frac{\log_2 e}{E[Q]} E[Q - E[Q]]^2 \text{ bits} . \quad (d)$$

The variance of Q is given by

$$\text{Var}[Q] = E[Q - E[Q]]^2 , \quad (e)$$

and thus we may rewrite eq. (d) as

$$E[N] \approx \log_2 E[Q] - \frac{\log_2 e}{E[Q]} \text{Var}[Q] \text{ bits} . \quad (f)$$

Since Q is geometrically distributed with a probability of success ε^{m-1} , the expectation and variance of Q is known to be

$$E[Q] = \varepsilon^{-(m-1)} \quad (g)$$

and

$$\text{Var}[Q] = \frac{1 - \varepsilon^{m-1}}{(\varepsilon^{m-1})^2} . \quad (h)$$

Plugging these values into eq. (f) yields

$$E[N] \approx \log_2 \varepsilon^{-(m-1)} - \frac{\log_2 e}{\varepsilon^{-(m-1)}} \frac{1 - \varepsilon^{(m-1)}}{(\varepsilon^{m-1})^2} \quad (i)$$

$$= -(m-1) \log_2 \varepsilon + \left(1 - \varepsilon^{-(m-1)}\right) \log_2 e \text{ bits} . \quad (j)$$

We are interested in the bits per element. There are m elements, so dividing by m results in

$$- \frac{m-1}{m} \log_2 \varepsilon + \frac{1 - \varepsilon^{-(m-1)}}{m} \text{ bits/element} . \quad (k)$$

Asymptotically, as $m \rightarrow \infty$, the expected bits per element goes to

$$- \log_2 \varepsilon . \quad (l)$$

□

By ??, ?? has an expected time complexity that grows exponentially as m grows The algorithm is intended to illustrate theoretical properties, not necessarily be used in practice.

6.2 Entropy

Proof.

$$\mathcal{H}(N) = - E[\log_2 p_N(N)] . \quad (a)$$

$$\mathcal{H}(N) = - \sum_{n=0}^{\infty} \log_2 p_N(n) p_N(n) \quad (b)$$

$$= - \sum_{n=0}^{\infty} \log_2 \left(q^{2^n-1} \left(1 - q^{2^n}\right) \right) p_N(n) . \quad (c)$$

$$\begin{aligned} \mathcal{H}(N) = & -\log_2 q \sum_{n=0}^{\infty} (2^n - 1) p_N(n) - \\ & \sum_{n=0}^{\infty} \log_2 (1 - q^{2^n}) p_N(n). \end{aligned} \quad (d)$$

$$\begin{aligned} \mathcal{H}(N) = & -\log_2 q \left[\sum_{n=0}^{\infty} (2^n p_N(n)) \right] - \\ & \sum_{n=0}^{\infty} \log_2 (1 - q^{2^n}) p_N(n). \end{aligned} \quad (e)$$

□

Theorem 10. *The random bit string that codes the Singular Hash Set is a maximum entropy coder for the approximate set abstract data type.*

Proof. The result immediately follows from the fact that the bit string is *incompressible*, and thus obtains maximum entropy. The hash h_k is the result of a random oracle and is thus uniformly distributed and, given that a bit string of length n is found, the probability that a particular b_n is found is uniformly distributed also. The only part of this that does not obtain maximum entropy is the particular bit length n of b_n . □

The only information contained in the encoding of the approximate set is given by the probability mass function of the random bit length N , which is a function of the cardinality of the finite set being approximated and the false positive rate ε .

Given a set \mathcal{S} , the random bit length N of bit string b_n has a probability mass concentrated around the theoretical lower-bound. Therefore, a *method-of-moments* estimator of the cardinality of the set \mathcal{S} is given by

$$|\hat{\mathcal{S}}| = -\frac{\text{BL}(\mathcal{S}^+)}{\log_2 \varepsilon}, \quad (45)$$

where BL is the bit length function and ε is the *false positive rate*.

We may trade space complexity for entropy if desired. For instance, if the policy is to search for a bit string of a length t much larger than the expected bit length, $-m \log_2 \varepsilon$, then

$$m < -\frac{t}{\log_2 \varepsilon}. \quad (46)$$

Thus, an estimator of the upper-bound on the cardinality is given by

$$\hat{m}_{\max} \approx -\frac{t}{\log_2 \varepsilon} \quad (47)$$

and an estimate of the *lower-bound* is 0 (the empty set) where the cardinality is uniformly distributed (maximum entropy) between 0 and \hat{m}_{\max} , which has an entropy given by

$$\log_2(1 + \hat{m}_{\max}). \quad (48)$$

The absolute space efficiency is now given by

$$\frac{m}{\hat{m}_{\max}}, \quad (49)$$

which is the *maximum* efficiency possible for an approximate set with a cardinality uniformly distributed between 0 and \hat{m}_{\max} .

Example 5 Suppose $t = rm, r > 1$, then

$$\hat{m}_{\max} \approx -\frac{rm}{\log_2 \varepsilon} \quad (\text{a})$$

which has an entropy given by

$$\log_2 \left(1 - \frac{rm}{\log_2 \varepsilon} \right) \approx \log_2 r + \log_2 m + \log_2 \frac{1}{\varepsilon} \quad (\text{b})$$

and an absolute space efficiency r .

Appendices

A Probability mass of random bit length

Algorithm 5: Bit length sampler of Singular Hash Set

in : m is the cardinality of the random set to approximate and ε is the false positive rate.

out : A *minimum bit length* n of a Singular Hash Set conditioned on a random set with cardinality m and a false positive rate ε .

```

1 function bit_length_sampler( $m, \varepsilon$ ):
2    $\mathcal{S} \leftarrow \emptyset$ ;
3   for  $i \leftarrow 1$  to  $m$  do
4      $x \leftarrow$  randomly draw a bit string from  $\mathcal{B}^*$  without replacement;
5      $\mathcal{S} \leftarrow \mathcal{S} \cup \{x\}$ ;
6   end
7    $(h_k, b_n) \leftarrow \text{make\_singular\_hash\_set}(\mathcal{S}, \varepsilon)$ ;
8   return  $k + n + \mathcal{O}(1)$ ;
```

For a particular n , each $b_n \in \mathcal{B}^n$ may fail to result in a perfect collision, therefore n is uncertain and takes on particular values with probabilities given by the following theorem.

Definition 13. The Singular Hash Set generator given by ?? finds a random string of a random bit length given by

$$N = \text{bit_length_sampler}(m, \varepsilon), \quad (50)$$

conditioned on a random set of cardinality m and a false positive rate ε .

Theorem 11. The random bit length N has a probability mass function given by

$$p_N(n \mid m, \varepsilon) = q^{2^n - 1} (1 - q^{2^n}), \quad (51)$$

where $q = 1 - \varepsilon^{m-1}$, m is the cardinality of the random set, and ε is the false positive rate.

Proof. Each iteration of the loop in ?? has a collision test which is Bernoulli distributed with a probability of success ε^{m-1} , where success denotes a perfect collision. We are interested in the random length N of the bit string when this outcome occurs.

For the random variable N to realize a value n , every bit string smaller than length n must fail and a bit string of length n must succeed. There are $2^n - 1$ bit strings smaller than length n and

each one fails with probability q , and so by the product rule the probability that they all fail is given by

$$q^{2^n-1} . \quad (\text{a})$$

Given that every bit string smaller than length n fails, what is the probability that every bit string of length n fails? There are 2^n bit strings of length n , each of which fails with probability q as before and thus by the product rule the probability that they all fail is q^{2^n} , whose complement, the probability that not all bit strings of length n fail, is given by

$$1 - q^{2^n} . \quad (\text{b})$$

By the product rule, the probability that every bit string smaller than length n fails and a bit string of length n succeeds is given by the product of (a) and (b),

$$q^{2^n-1} (1 - q^{2^n}) . \quad (\text{c})$$

For Equation (c) to be a probability mass function, two conditions must be met. First, its range must be a subset of $[0, 1]$. Second, the summation over its domain must be 1.

The first case is trivially shown by the observation that q is a positive number between 0 and 1 and therefore any non-negative power of q is positive number between 0 and 1.

The second case is shown by calculating the infinite series

$$S = \sum_{n=0}^{\infty} q^{2^n-1} (1 - q^{2^n}) \quad (\text{d})$$

$$= \sum_{n=0}^{\infty} q^{2^n-1} - q^{2^{n+1}-1} . \quad (\text{e})$$

Explicitly evaluating this series for the first 4 terms reveals a telescoping sum given by

$$S = (1 - q) + (q - q^3) + (q^3 - q^7) + (q^7 - q^{15}) + \dots , \quad (\text{f})$$

where everything cancels except 1.

Note that a simpler proof is given by

$$N = \log_2 Q . \quad (\text{g})$$

Thus, the probability mass function is given by

$$p_N(n \mid m, \varepsilon) = P[N = n] \quad (\text{h})$$

$$= P[\log_2 Q = n] \quad (\text{i})$$

$$= P[Q = 2^n] \quad (\text{j})$$

$$= p_Q(2^n \mid m, \varepsilon) \quad (\text{k})$$

$$= \varepsilon^{m-1} (1 - \varepsilon^{m-1})^{2^n} \quad (\text{l})$$

□

References

- [1] Alexander Towell. The approximate set abstract data type. . URL ?
- [2] Alexander Towell. The oblivious object type. . URL ?