

An algebra of random approximate sets

with derivations of higher-order random approximate sets induced by set-theoretic operations on random approximate sets with corresponding random binary classification measures.

Alexander Towell
atowell@siue.edu

Abstract

We define a *random approximate set* model and the probability space that follows. A random approximate set is a *probabilistic* set generated to *approximate* another set of objective interest. We derive several properties that follow from this definition, such as the expected *precision* in information retrieval. Finally, we demonstrate an application of approximate sets, approximate Encrypted Search with queries as a Boolean algebra, which generates random approximate result sets.

Contents

1	Introduction	1
2	Algebra of sets	1
2.1	Boolean algebras	2
3	Random approximate set model	3
3.1	First-order model	5
3.1.1	Probability space	5
3.2	Higher-order model	6
4	Probability distributions of parameters for <i>first-order</i> approximations	7
4.1	Asymptotic limits	9
5	Probability distributions of parameters for <i>higher-order</i> approximations	10
5.1	Compositions of random approximate sets	10
6	Distribution of binary classification measures	12
7	Uncertain rate distortions	15
7.1	First-order model	16
8	Data types that model random approximate sets	19
8.1	Deterministic value constructors	19
8.2	Space complexity	21
8.2.1	Space efficiency of <i>unions</i> and <i>differences</i>	22
8.3	code tmp	23
8.4	Algebra of sets	23
9	Application: approximating Boolean search	24
9.1	Random approximate Boolean search	25
A	Proof of corollary 4.2.2	25
B	Proof of theorem 6.1	26
C	Proof of theorem 5.3	27
D	Sampling distribution of arbitrary functions	29

1 Introduction

An *approximate set* is a set that approximates another set of objective interest. It is *approximate* because with respect to the objective set, there are two types of errors, *false positives* and *false negatives*. The *Bloom filter* is a popular example of a data structure that models *positive random approximate sets* with *false positives* due to *rate distortion*.

In section 2, we define the algebra of sets.

In section 3, we provide a formal definition of the *random approximate set* model, in which the false positive and false negative rates are *expectations*. We describe the axioms of the random approximate set model such that, if satisfied, also satisfy the axioms of the approximate algebra of sets. We further derive the probability distribution of random approximate sets entailed by the axioms.

In section 4, we derive the rate-distortion random variables that are fundamental to the approximate set model, such as the false positive rate.

In ??, we provide a detailed treatment on distributions that are induced by functions that depend on random approximate sets, e.g., in section 5.1 we derive the probability distribution of random approximate sets that are generated from arbitrary set-theoretic operations on random approximate sets and in section 6 we derive several well-known binary classification performance measures of random approximate sets as a function of their error rates, such as *positive predictive value*.

In section 7, we provide the probabilistic model for random approximate sets with *uncertain* rate distortions, such as an uncertain false positive rate.

In section 8, we provide a treatment on the random approximate set model as an abstract data type and show how that, if the generative algorithm of an approximate set model is deterministic, the random approximate set model quantifies our ignorance or uncertainty.

Finally, in section 9, we consider Encrypted Search with secure indexes based on random approximate sets. To prove various properties of this model, such as expected precision, we only need to show that the *result sets* are approximate sets of the *objective* results and all the results immediately follow.

2 Algebra of sets

A *set* is an unordered collection of distinct elements. If we know the elements in a set, we may denote the set by these elements, e.g., $\{a, c, b\}$ denotes a set whose members are exactly a , b , and c .

Two sets of particular importance are the empty set, denoted by \emptyset , which has no members, and the *universal set*, in which every element of interest is a member.

A *finite set* has a finite number of elements. For example, $\{1, 3, 5\}$ is a finite set with three elements. When sets \mathcal{A} and \mathcal{B} are *isomorphic*, denoted by $\mathcal{A} \cong \mathcal{B}$, they can be put into a one-to-one correspondence (bijection), e.g., $\{b, a, c\} \cong \{1, 2, 3\}$.

The cardinality of a finite set \mathcal{A} is the number of elements in the set, denoted by $|\mathcal{A}|$, e.g., $|\{1, 3, 5\}| = 3$. A *countably infinite set* is isomorphic to the set of *natural numbers* $\mathbb{N} := \{1, 2, 3, 4, 5, \dots\}$.

Given two elements a and b , an ordered pair of a then b is denoted by $\langle a, b \rangle$, where $\langle a, b \rangle = \langle c, d \rangle$ if and only if $a = c$ and $b = d$. Ordered pairs are non-commutative and non-associative, i.e., $\langle a, b \rangle \neq \langle b, a \rangle$ if $a \neq b$ and $\langle a, \langle b, c \rangle \rangle \neq \langle \langle b, a \rangle, c \rangle$.

Related to the ordered pair is the Cartesian product.

Definition 2.1. *The set $\mathcal{X} \times \mathcal{Y} := \{\langle x, y \rangle : x \in \mathcal{X} \wedge y \in \mathcal{Y}\}$ is the Cartesian product of sets \mathcal{X} and \mathcal{Y} .*

By the non-commutative and non-associative property of ordered pairs, the Cartesian product is non-commutative and non-associative. However, they are isomorphic, i.e., $\mathcal{X} \times \mathcal{Y} \cong \mathcal{Y} \times \mathcal{X}$.

A *tuple* is a generalization of order pairs which can consist of an arbitrary number of elements, e.g., $\langle x_1, x_2, \dots, x_n \rangle$.

Definition 2.2 (*n*-fold Cartesian product). *The n-ary Cartesian product of sets $\mathcal{X}_1, \dots, \mathcal{X}_n$, is given by $\mathcal{X}_1 \times \dots \times \mathcal{X}_n = \{\langle x_1, \dots, x_n \rangle : x_1 \in \mathcal{X}_1 \wedge \dots \wedge x_n \in \mathcal{X}_n\}$.*

Note that $\mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{X}_3 \cong \mathcal{X}_1 \times (\mathcal{X}_2 \times \mathcal{X}_3) \cong (\mathcal{X}_1 \times \mathcal{X}_2) \times \mathcal{X}_3$, thus we may implicitly convert between them without ambiguity.

If each set in the *n*-ary Cartesian product is the same, the power notation may be used, e.g., $\mathcal{X}^3 := \mathcal{X} \times \mathcal{X} \times \mathcal{X}$. As special cases, $\mathcal{X}^0 := \{\emptyset\}$ and $\mathcal{X}^1 := \mathcal{X}$.

A *binary relation* over sets \mathcal{A} and \mathcal{B} is any subset of $\mathcal{A} \times \mathcal{B}$. A fundamental relation is the member-of relation, where $x \in \mathcal{A}$ denotes that an object x is a member of a set \mathcal{A} . A set \mathcal{A} is a *subset* of a set \mathcal{B} if every member of \mathcal{A} is a member \mathcal{B} , denoted by $\mathcal{A} \subseteq \mathcal{B}$. The subset relation forms a *partial order*, i.e., if $\mathcal{A} \subseteq \mathcal{B}$ and $\mathcal{B} \subseteq \mathcal{C}$ then $\mathcal{A} \subseteq \mathcal{C}$ and if $\mathcal{A} \subseteq \mathcal{B}$ and $\mathcal{B} \subseteq \mathcal{A}$ then \mathcal{A} and \mathcal{B} are *equal*, denoted by $\mathcal{A} = \mathcal{B}$.

Definition 2.3. *Set builder notation*

Definition 2.4. *A function of type $\mathcal{X} \mapsto \mathcal{Y}$ is a binary relation on $\mathcal{X} \times \mathcal{Y}$ with the constraint that each $x \in \mathcal{X}$ is paired with exactly one $y \in \mathcal{Y}$.*

A function of type $\mathcal{X} \mapsto \mathcal{Y}$ has a domain \mathcal{X} and a codomain \mathcal{Y} . Since every $x \in \mathcal{X}$, given a pair $\langle x, y \rangle \in f$, y may also be denoted by $f(x)$.

The *power set* of a set \mathcal{A} , denoted by $2^{\mathcal{A}}$, is the set of sets that contains all of the possible subsets of \mathcal{A} , e.g., $2^{\{a,b\}} = \{\emptyset, \{a\}, \{b\}, \{a, b\}\}$.

A predicate is a function that maps elements in its domain to true (denoted by 1) or false (denoted by 0). A predicate function of particular importance is the indicator function

$$\mathbb{1}_{\mathcal{A}}: \mathcal{X} \mapsto \{0, 1\} \quad (2.1)$$

defined as

$$\mathbb{1}_{\mathcal{A}}(x) := \begin{cases} 0 & \text{if } x \notin \mathcal{A}, \\ 1 & \text{if } x \in \mathcal{A}. \end{cases} \quad (2.2)$$

The indicator function admits the construction of predicates for any relation, e.g., a binary predicate P for a binary relation $\mathcal{R} \subseteq \mathcal{A} \times \mathcal{B}$ is defined as $P(x_1, x_2) := \mathbb{1}_{\mathcal{R}}(\langle x_1, x_2 \rangle)$. Denoting the *universal set* by \mathcal{X} , all the relations mentioned previously are *binary predicates*, such as $\in: \mathcal{X} \times 2^{\mathcal{X}} \mapsto \{0, 1\}$ and $\subseteq: 2^{\mathcal{X}} \times 2^{\mathcal{X}} \mapsto \{0, 1\}$.

Some important operations on sets are described next. The *union* operator, $\cup: 2^{\mathcal{X}} \times 2^{\mathcal{X}} \mapsto 2^{\mathcal{X}}$, is defined as

$$\mathcal{A} \cup \mathcal{B} := \{x \in \mathcal{X} \mid x \in \mathcal{A} \vee x \in \mathcal{B}\} \quad (2.3)$$

where \vee is the logical-connective *or*. The *intersection* operator, $\cap: 2^{\mathcal{X}} \times 2^{\mathcal{X}} \mapsto 2^{\mathcal{X}}$, is defined as

$$\mathcal{A} \cap \mathcal{B} := \{x \in \mathcal{X} \mid x \in \mathcal{A} \wedge x \in \mathcal{B}\} \quad (2.4)$$

where \wedge is the logical-connective *and*. If $\mathcal{A} \cap \mathcal{B} = \emptyset$, then we say \mathcal{A} and \mathcal{B} are *disjoint* sets.

The *relative complement* (set-difference) operator, $\setminus: 2^{\mathcal{X}} \times 2^{\mathcal{X}} \mapsto 2^{\mathcal{X}}$, is defined as

$$\mathcal{A} \setminus \mathcal{B} := \{x \in \mathcal{X} \mid x \in \mathcal{A} \wedge x \notin \mathcal{B}\}. \quad (2.5)$$

The relative complement $\mathcal{X} \setminus \mathcal{A}$ is denoted by $\overline{\mathcal{A}}$ and is called the *complement* of \mathcal{A} .

2.1 Boolean algebras

An *algebra* denotes a mathematical structure in which a certain set of axioms hold. A *Boolean algebra* is given by the following definition.

Definition 2.5. A *Boolean algebra* is a six-tuple $(\mathcal{A}, \wedge, \vee, \neg, 0, 1)$ where \mathcal{A} is a set, \wedge is the binary meet operation, \vee is the binary join operation, \neg is the unary complement operation, 0 is the bottom element, and 1 is the top element such that $\forall a, b, c \in \mathcal{A}$ the following axioms hold:

1. *Associativity:* $a \vee (b \vee c) = (a \vee b) \vee c$ and $a \wedge (b \wedge c) = (a \wedge b) \wedge c$.
2. *Commutativity:* $a \vee b = b \vee a$ and $a \wedge b = b \wedge a$.
3. *Identity:* $a \vee 0 = a$ and $a \wedge 1 = a$.
4. *Distributivity:* $a \vee (b \wedge c) = (a \vee b) \wedge (a \vee c)$ and $a \wedge (b \vee c) = (a \wedge b) \vee (a \wedge c)$.
5. *Complementation:* $a \vee \neg a = 1$ and $a \wedge \neg a = 0$.

Every valid proposition in a Boolean algebra is derivable from the axioms in definition 2.5. A particularly useful result is *De Morgan's laws*,

$$a \vee b = \neg(\neg a \wedge \neg b) \quad (2.6)$$

and

$$a \wedge b = \neg(\neg a \vee \neg b). \quad (2.7)$$

Postulate 2.1. Given the universal set \mathcal{U} and a set $\Sigma \subseteq 2^{\mathcal{U}}$ that is closed under unions, intersections, and complements, $(\Sigma, \cup, \cap, \overline{\cdot}, \emptyset, \mathcal{U})$ is a Boolean algebra.

Trivially, $\Sigma = 2^{\mathcal{U}}$ forms a Boolean algebra, but later we demonstrate that implementations of the *random approximate set* model may form a Boolean algebra over some closed subset $\Sigma \subset 2^{\mathcal{U}}$.

The algebra of bit-wise operations on vectors of u bits is given by $(\{0, 1\}^u, \wedge, \vee, \neg, \mathbf{0}, \mathbf{1})$ where \wedge is bit-wise *and*, \vee is bit-wise *or*, \neg is bit-wise *negation*, $\mathbf{0}$ is vector of all zeros, and $\mathbf{1}$ is vector of all ones.

A bijection between the algebra of sets and the algebra of bit vectors is given by the following definition.

Definition 2.6. Suppose there is some total order on \mathcal{U} , $u = |\mathcal{U}|$, such that the j -th ranked element may be denoted by $x_{(j)}$. A bijection between the Boolean algebras $(2^{\mathcal{U}}, \cap, \cup, \overline{\cdot}, \emptyset, \mathcal{U})$ and $(\{0, 1\}^u, \wedge, \vee, \neg, \mathbf{0}, \mathbf{1})$ is given by mapping $\mathcal{X} \in 2^{\mathcal{U}}$ to $\mathbf{a} \in \{0, 1\}^u$ where $a_j = \mathbb{1}_{\mathcal{X}}(x_{(j)})$. Additionally, $\vee \leftrightarrow \cup$, $\wedge \leftrightarrow \cap$, $\neg \leftrightarrow \overline{\cdot}$, $\mathbf{0} \leftrightarrow \emptyset$, and $\mathbf{1} \leftrightarrow \mathcal{U}$.

This bijection allows us to use either representation interchangeably.

3 Random approximate set model

The concept of a random approximate set depends upon the concept of an *approximate set*.

Given an objective set \mathcal{S} , any element that is a member of \mathcal{S} is denoted a *positive* of \mathcal{S} and any element that is *not* a member of \mathcal{S} is denoted a *negative* of \mathcal{S} .

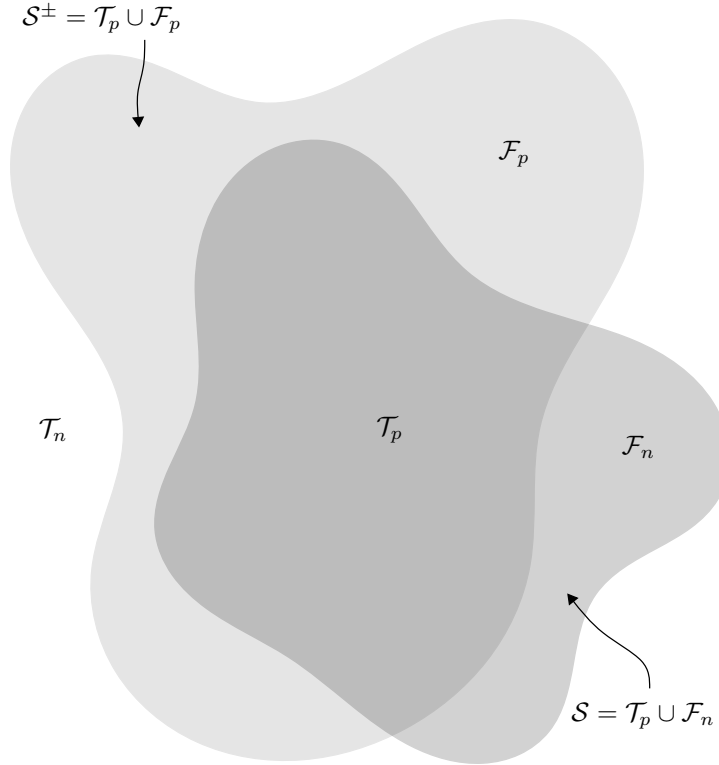
A set that is used as an *approximation* of \mathcal{S} may be denoted by \mathcal{S}^\pm . If the *only* information we have about \mathcal{S} is given by \mathcal{S}^\pm , then we may perform membership tests on \mathcal{S}^\pm to *predict* the members (or non-members) of \mathcal{S} .

There are two ways a binary prediction can be false.

1. A *false positive* occurs if a negative of the objective set is predicted to be a positive. False positives are also known as *type I errors*. The complement of false positives are *true negatives*.
2. A *false negative* occurs if a positive of the objective set is predicted to be a negative. False negatives are also known as *type II errors*. The complement of false negatives are *true positives*.

Suppose we have an objective set \mathcal{S} and an approximation \mathcal{S}^\pm . If we denote the set of false positives by \mathcal{F}_p , true positives by \mathcal{T}_p , false negatives by \mathcal{F}_n , and true negatives by \mathcal{T}_n , then the objective set \mathcal{S} is equal to $\mathcal{F}_n \cup \mathcal{T}_p$ and the approximate set \mathcal{S}^\pm is equal to $\mathcal{T}_p \cup \mathcal{F}_p$. See fig. 1 for an illustration.

Figure 1: An approximate set \mathcal{S}^\pm of an objective set \mathcal{S}



If we only have access to the approximation \mathcal{S}^\pm , we cannot partition the universe into the sets \mathcal{F}_p , \mathcal{T}_p , \mathcal{F}_n , and \mathcal{T}_n as demonstrated in fig. 1. However, we can quantify the degree of *uncertainty* about the elements that are predicted to be positive or negative.

The false positive and true negative *rates* are given by the following.

Definition 3.1. The false positive rate is the proportion of predictions that are false positives as given by

$$\hat{\varepsilon} = \frac{f_p}{f_p + t_n}, \quad (3.1)$$

where f_p is the number of false positives and t_n is the number of true negatives. In a complementary manner, the true negative rate is $\hat{\eta} = 1 - \hat{\varepsilon}$.

The true positive and false negative *rates* are given by the following.

Definition 3.2. The true positive rate is the proportion of predictions that are true positives as given by

$$\hat{\tau} = \frac{t_p}{t_p + f_n}, \quad (3.2)$$

where f_n is the number of false negatives and t_p is the number of true positives. In a complementary manner, the false negative rate is $\hat{\omega} = 1 - \hat{\tau}$.

The *probabilities* of the four possible predictive outcomes are given by table 1.

	positive	negative
predict positive	$\hat{\tau} = 1 - \hat{\omega}$	$\hat{\varepsilon} = 1 - \hat{\eta}$
predict negative	$\hat{\omega} = 1 - \hat{\tau}$	$\hat{\eta} = 1 - \hat{\varepsilon}$

Table 1: The 2×2 contingency table of outcomes for approximate sets.

In the *random* approximate set model, we do not describe any particular approximation, but rather we describe the statistical properties of processes that *generate* approximations.

The random false positive rate \mathcal{E} and false negative rate \mathcal{W} have supports in the Borel set of $[0, 1]$.

The *zero-th order* generative model for sets is not generally known, but we denote the zero-th order model by R . We denote the *first-order* random approximate set generative model by R^\pm . The joint distribution of R^\pm , \mathcal{E} , \mathcal{W} , and R given a universal set \mathcal{U} has a probability density

$$f_{R^\pm, \mathcal{E}, \mathcal{W}, R}(\mathcal{Y}, \varepsilon, \omega, \mathcal{X} | \mathcal{U}). \quad (3.3)$$

By the axioms of probability theory, this may be decomposed into

$$f_{R^\pm, \mathcal{E}, \mathcal{W}, R}(\mathcal{Y}, \varepsilon, \omega, \mathcal{X} | \mathcal{U}) = f_{R^\pm | \mathcal{E}, \mathcal{W}, R}(\mathcal{Y} | \varepsilon, \omega, \mathcal{X} | \mathcal{U}) f_{\mathcal{E}, \mathcal{W} | R}(\varepsilon, \omega | \mathcal{X}) f_R(\mathcal{X} | \mathcal{U}). \quad (3.4)$$

We typically omit the explicit reference to \mathcal{U} , since it may usually be understood as implicit to the model.

The object of central interest is the distribution of R^\pm given R . The conditional distribution of R^\pm given $R = \mathcal{X}$ is denoted by \mathcal{X}^\pm . By the axioms of probability,

$$f_{\mathcal{X}^\pm, \mathcal{E}, \mathcal{W}}(\mathcal{Y}, \varepsilon, \omega) = f_{\mathcal{X}^\pm | \mathcal{E}, \mathcal{W}}(\mathcal{Y} | \varepsilon, \omega) f_{\mathcal{E}, \mathcal{W} | R}(\varepsilon, \omega | \mathcal{X}). \quad (3.5)$$

The random false positive and false negative rates conditioned on $R = \mathcal{X}$ are respectively given by

$$\mathcal{E} = \frac{1}{|\mathcal{X}|} \sum_{x \in \bar{\mathcal{X}}} \mathbb{1}_{\mathcal{X}^\pm}(x) \quad (3.6)$$

and

$$\mathcal{W} = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \mathbb{1}_{\mathcal{X}^\pm}(x). \quad (3.7)$$

\mathcal{A}^\pm conditioned on $\mathcal{E} = a$ and $\mathcal{W} = b$ is a random approximate set with the indicated false positive and false negative rates. If the rates happen to pick out a specific set in the support, then the result is a degenerate distribution, e.g., \mathcal{A}^\pm given $\mathcal{E} = 0$ and $\mathcal{W} = 0$ is degenerate where all probability mass is assigned to \mathcal{A} .

We denote the distributions of \mathcal{X}^\pm given $E[\mathcal{E}] = \varepsilon$ and \mathcal{X}^\pm given $E[\mathcal{W}] = \omega$ respectively by $\mathcal{X}_\varepsilon^-$ and \mathcal{X}_ω^+ . An object of central interest is the distribution of \mathcal{X}^\pm given $E[\mathcal{E}] = \varepsilon$ and $E[\mathcal{W}] = \omega$, denoted by

$$\mathcal{X}_\varepsilon^\omega. \quad (3.8)$$

If we *sample* from $\mathcal{A}_\varepsilon^\omega$, some set $\mathcal{Y} \in 2^{\mathcal{U}}$ with false positive rate a and false negative rate b will be realized with probability $f_{\mathcal{A}_\varepsilon^\omega | \mathcal{E}, \mathcal{W}}(\mathcal{Y} | a, b)$. However, as the number of samples goes to infinity, the mean false positive and false negative rates go to ε and ω respectively.

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

Definition 3.3. A random approximate set $\mathcal{A}_\varepsilon^-$ is a random positive approximate set denoted by \mathcal{A}^+ .

By this definition, any instance of \mathcal{A}^+ is a *superset* of \mathcal{A} .

As shown in ??, positive approximate sets are closed under unions and intersections but not complements. We introduce the *negative approximate set* as a natural consequence.

Definition 3.4. A random approximate set \mathcal{A}_ω^+ is a random negative approximate set denoted by \mathcal{A}^- .

By this definition, any instance of \mathcal{A}^- is a *subset* of \mathcal{A} .

Negative approximate sets are *closed* under unions and intersections but not complements. The complement of a random positive (negative) approximate set is a random negative (positive) approximate set.

3.1 First-order model

The first-order random approximate set model makes the following assumption about the joint distribution of \mathcal{E} , \mathcal{W} , and \mathbb{R} .

Axiom 1. *The random variables \mathcal{E} and \mathcal{W} are conditionally independent given \mathbb{R} .*

By axiom 1 and by the axioms of probability,

$$f_{\mathcal{X}^\pm, \mathcal{E}, \mathcal{W}}(\mathcal{Y}, \varepsilon, \omega) = f_{\mathcal{X}^\pm | \mathcal{E}, \mathcal{W}}(\mathcal{Y} | \varepsilon, \omega) f_{\mathcal{E} | \mathbb{R}}(\varepsilon | \mathcal{X}) f_{\mathcal{W} | \mathbb{R}}(\omega | \mathcal{X}). \quad (3.9)$$

The following axioms complete the probabilistic model for *first-order* random approximate sets.

Axiom 2. *The outcome of a membership test on any element in the negative set is an independent and identically distributed Bernoulli trial with a mean ε ,*

$$\mathbb{P}[\mathbb{1}_{\mathcal{A}_\tau^\varepsilon}(x) | \neg \mathbb{1}_{\mathcal{A}}(x)] = \varepsilon. \quad (3.10)$$

Axiom 3. *The outcome of a membership test on any element in the negative set is an independent and identically distributed Bernoulli trial with a mean τ ,*

$$\mathbb{P}[\mathbb{1}_{\mathcal{A}_\tau^\varepsilon}(x) | \mathbb{1}_{\mathcal{A}}(x)] = \tau. \quad (3.11)$$

It may not be possible, desirable, or practical to observe these rates, e.g., the objective set may not be knowable from the given information. In section 4 we derive the probability distributions for characteristics like the false positive rate. Thus, for instance, we may provide a *confidence interval* which contains the false positive rate with some probability α which is a function of parameters like the expected false positive rate ε .

Every statistical property of the random approximate set model (first-order and higher-order) is entailed by axioms 2 and 3. Furthermore, these assumptions generally hold in practice, e.g., the Bloom filter[?] and Perfect hash filter[?] are two separate implementations¹ of the random positive approximate set in which these assumptions hold.

Suppose the first-order random approximate sets are over the universal set \mathcal{U} . Then, over the Boolean algebra $(2^{\mathcal{U}}, \cup, \cap, \bar{\cdot}, \emptyset, \mathcal{U})$, the approximate sets formed are no longer *first-order* approximations. In ??, we describe such *higher-order* models.

3.1.1 Probability space

Suppose the universal set is \mathcal{U} and we have some process that generates approximations of some objective set \mathcal{A} that is compatible with the axioms of the random approximate set model.

The process generates subsets of \mathcal{U} , or alternatively, the *sample space* is $\Sigma = 2^{\mathcal{U}}$. A primary objective in *probability modeling* is assigning *probabilities to events*. Suppose we have some *probability function* $\mathbb{P}: \Sigma \mapsto [0, 1]$. The *probability* of some event $\mathcal{A} \in \Sigma$ is denoted by $\mathbb{P}[\mathcal{A}]$.

These are the *elementary events* of the probability space. The random approximate set model given $\mathbb{R} = \mathcal{Y}$ is given by the *probability space*

$$(\Omega = 2^{\mathcal{U}}, 2^\Omega, \mathbb{P}), \quad (3.12)$$

where Ω is the *sample space*, 2^Ω is the set of all events, and $\mathbb{P}: 2^\Omega \mapsto [0, 1]$ is the probability set function.

Consider an objective set \mathcal{A} and a random approximate set and suppose we are uncertain about which elements are their respective members. We model the uncertainty of the elements of \mathcal{A} by the Boolean random vector $\mathbf{A} = \langle A_1, \dots, A_u \rangle$ where $A_j = \mathbb{1}_{\mathcal{A}}(x_{(j)})$ for $j = 1, \dots, u$. Similarly, we model the uncertainty of the elements of \mathcal{A}^\pm by $\mathbf{A}_\tau^\varepsilon = \langle A_1^\pm, \dots, A_u^\pm \rangle$.

The joint probability that $\mathbf{A}_\tau^\varepsilon = \mathbf{x}$ and $\mathbf{A} = \mathbf{y}$ is denoted by $\mathbb{P}[\mathbf{A}_\tau^\varepsilon = \mathbf{x}, \mathbf{A} = \mathbf{y}]$. By the axioms of probability, the joint probability may be rewritten as

$$\mathbb{P}[\mathbf{A}_\tau^\varepsilon = \mathbf{x}, \mathbf{A} = \mathbf{y}] = \mathbb{P}[\mathbf{A}_\tau^\varepsilon = \mathbf{x} | \mathbf{A} = \mathbf{y}] \mathbb{P}[\mathbf{A} = \mathbf{y}]. \quad (3.13)$$

By axioms 2 and 3, A_j^\pm is only dependent on A_j for $j = 1, \dots, u$ and thus by the axioms of probability

$$\mathbb{P}[\mathbf{A}_\tau^\varepsilon = \mathbf{x}, \mathbf{A} = \mathbf{y}] = \mathbb{P}[\mathbf{A} = \mathbf{y}] \prod_{j=1}^u \mathbb{P}[A_j^\pm = x_j | A_j = y_j]. \quad (3.14)$$

If it is given that $\mathbf{A}_\tau^\varepsilon = \mathbf{y}$, i.e., the elements in the objective set are known, by the axioms of probability the conditional probability is

$$\mathbb{P}[\mathbf{A}_\tau^\varepsilon = \mathbf{x} | \mathbf{A} = \mathbf{y}] = \prod_{j=1}^u \mathbb{P}[A_j^\pm = x_j | A_j = y_j] \quad (3.15)$$

where $\varepsilon = \mathbb{P}[A_j^\pm = 1 | A_j = 0]$ and $\tau = \mathbb{P}[A_j^\pm = 1 | A_j = 1]$.

¹There may be a difference in that the algorithm may be deterministic; we address this point in ??.

The relative frequency of any event \mathbf{x} in $\{0,1\}^u$ converges to $P[\mathbf{X}^\pm = \mathbf{x} \mid \mathbf{y}^\pm]$ as the number of times the random approximate set of \mathbf{y}^\pm is generated goes to infinity.

Consider the following example.

Example 1 Suppose the universal set is $\{x_1, x_2\}$ and consider the distribution of the first-order random approximate set $\{x_1\}_\varepsilon^\omega$. The probability mass function $p_{\{x_1\}_\varepsilon^\omega}$ is given by

$$P_{\{x_1\}_\varepsilon^\omega}(\mathcal{X}) = \begin{cases} \omega(1 - \varepsilon) & \mathcal{X} = \emptyset, \\ \omega\varepsilon & \mathcal{X} = \{x_2\}, \\ (1 - \omega)(1 - \varepsilon) & \mathcal{X} = \{x_1\}, \\ (1 - \omega)\varepsilon & \mathcal{X} = \{x_1, x_2\}. \end{cases} \quad (\text{a})$$

3.2 Higher-order model

Composing *random approximate sets* over the Boolean algebra $(\Sigma, \cup, \cap, \bar{\cdot}, \emptyset, \mathcal{U})$, where $\Sigma \subseteq 2^{\mathcal{U}}$ since, for instance, if a *deterministic* algorithm implements the model some elements in $2^{\mathcal{U}}$ may not be reachable. As a result, to satisfy the identity and complementation axioms required by Boolean algebras, we make \emptyset and \mathcal{U} available in the model as special cases.

Remark. Alternatively, these axioms may be satisfied by making the empty set and the universal set degenerate cases, i.e., $P[\emptyset^\pm = \emptyset] = 1$ and $P[\mathcal{U}^\pm = \mathcal{U}] = 1$. \triangle

Furthermore, we may replace any of the operators in the Boolean algebra with *random approximations* that model the noisy or rate-distorted channel previously described, i.e., these operators may themselves be constructors for random approximate sets, e.g., $\mathcal{A} \cup_\varepsilon^\tau \mathcal{B} \sim (\mathcal{A} \cup \mathcal{B})_\tau^\varepsilon$ where \cup_ε^τ maps negatives to positives with probability ε and maps positives to negatives with probability τ .

Given two sets \mathcal{X} and \mathcal{Y} , the set of all possible functions from domain \mathcal{X} to codomain \mathcal{Y} is denoted by $\mathcal{X} \mapsto \mathcal{Y}$ (or $\mathcal{X}^{\mathcal{Y}}$ since there are a total of $|\mathcal{X}|^{|\mathcal{Y}|}$ functions in the set).²

A particular function in the set $\mathcal{X} \mapsto \mathcal{Y}$ may be given a label f and we denote that it is a function in this set with the notation $f: \mathcal{X} \mapsto \mathcal{Y}$.

A natural mapping is provided by the *identity* function $\text{id}: 2^{\mathcal{X}} \mapsto 2^{\mathcal{X}}$, which is defined as

$$\text{id}(\mathcal{A}) := \mathcal{A}. \quad (3.16)$$

However, suppose we only have an *approximation* of the identity function, denoted by $\text{id}_\varepsilon^\tau$, such that $\text{id}_\varepsilon^\tau(\mathcal{A}) \sim \mathcal{A}_\varepsilon^\tau$. Then $\text{id}_\varepsilon^\tau$ generates sets consistent with the random approximate set model.

If we compose random approximate sets, then we have *higher-order* random approximate sets.

Theorem 3.1. The composition of random approximate identity functions $\text{id}_\varepsilon^\tau \circ \text{id}_{\varepsilon'}^{\tau'}$ generates random approximate sets with a true positive rate $\tau\tau' + \omega\varepsilon'$ and false positive rate $\varepsilon\tau' + \eta\varepsilon'$.

Definition 3.5. The iterated function f^k is defined as k compositions of f where f^0 denotes the (non-random) identity function.

The composition $(\text{id}_\varepsilon^\tau)^k$ generates k -th order random approximate sets where the *zero-th* order random approximation is the identity, i.e., $(\text{id}^\pm)^0 = \text{id}$.

The function being approximately may take other forms, like *set-complementation* or *set-union*, e.g., let $\cup: 2^{\mathcal{X}} \times 2^{\mathcal{X}} \mapsto 2^{\mathcal{X}}$ be approximated by $\cup_\varepsilon^\tau: 2^{\mathcal{X}} \times 2^{\mathcal{X}} \mapsto 2^{\mathcal{X}}$. Then, $\mathcal{A} \cup_\varepsilon^\tau \mathcal{B}$ is a random approximate set of $\mathcal{A} \cup \mathcal{B}$ as before. However, $\cup_\varepsilon^\tau \circ \text{id}_\varepsilon^\tau$ generates second-order random approximate sets.

Suppose we have an iterable set that is the output of some random approximation of some objective set of interest. We may wish to apply a more space-efficient data structure for random approximate sets, such as a Bloom filter[?]. In this case, the result is a *second-order* random approximate set; that is, a random approximate set of a random approximate set.

Theorem 3.2. Given a random approximate set $\mathcal{A}_{\varepsilon_1}^{\tau_1}$, a random approximation of $\mathcal{A}_{\varepsilon_1}^{\tau_1}$ with a false positive rate ε_2 and true positive rate τ_2 is a second-order random approximate set of \mathcal{A} with a false positive rate $\varepsilon = \varepsilon_1\tau_2 + \eta_1\varepsilon_2$ and true positive rate $\tau = \tau_1\tau_2 + \omega_1\varepsilon_2$, denoted by $\mathcal{A}^{\sigma^2}(\tau, \varepsilon)$.

This result may be recursively applied to derive arbitrary k -th order random approximate sets as given by

$$\mathcal{A}^{\sigma^k} = \left(\mathcal{A}^{\sigma^{k-1}}\right)^\sigma \quad (3.17)$$

where the zero-th order $\mathcal{A}^{\sigma^0} = \mathcal{A}$.

²The domain \mathcal{X} may be a Cartesian product of other sets, e.g., $\mathcal{X}_1 \times \mathcal{X}_2 \mapsto \mathcal{Y}$ denotes a set of binary functions.

In an *algebra of set* $(2^{\mathcal{U}}, \cap, \cup, \bar{\cdot}, \mathcal{U}, \emptyset)$, we may compose sets to form new sets. When these sets model *random approximate sets*, then their compositions model *higher-order random approximate sets*, e.g., $\mathcal{A}_{\varepsilon_1}^{\tau_1} \cup \mathcal{B}_{\varepsilon_2}^{\tau_2}$ models a higher-order random approximate set which does not obey the model described in ??; rather, it partitions the negative set such that each partition may have a different false positive rate and similarly for the positive set.

TODO: not that while each partition may have a different false positive rate, it obeys the first-order model of that partition. this is the basis of the zeroth order -> first order -> higher order model, where higher orders (after first) are a result of the Boolean algebra, i.e., unions, intersections, complements, etc. TODO: let's pull the set-complement into the first order model, since if given a first order RAS, then its complement is also a first-order RAS. TODO: get some venn diagrams in this. have one for the first-order model, now let's extend it to, say, a union of two approximate sets.

$$\cup^{\pm} \mathcal{A}_{\varepsilon_1}^{\tau_1} [\mathcal{B}_{\varepsilon_2}^{\tau_2}]$$

This contrasts with $(\mathcal{A} \cup \mathcal{B})_{\varepsilon}^{\tau}$, in which the false positive rate for the negative elements are uniformly distributed and likewise for the positive elements. We call these random approximate sets *first-order approximations*. For completeness, the *zeroth-order* are the objective sets, e.g., the zeroth-order approximation of $\mathcal{A} \cup \mathcal{B}$ is $\mathcal{A} \cup \mathcal{B}$. The complexity of the probabilistic model increases as the order increases.

If we have a data structure that models random approximate sets and if we transmit an objective set over a noisy channel in which positives become negatives with some probability and likewise for negatives, then any constructed random approximate set will be a higher-order random approximate set.

TODO: probabilistic model of unions and complements of random approximate sets. Grab the stuff from the parameter distribution section. It just partitions the negative and positive sets so that the predictive test for negatives are no longer uniformly distributed and likewise for positives.

4 Probability distributions of parameters for *first-order approximations*

The (first-order) random approximate sets are *parameterized* by the *expected* rates of two types of error, false negative and false positive rates. In this section, we derive the distribution for these rates.

A random variable $W: \Sigma \mapsto \mathcal{Y}$ is a function that maps outcomes in the σ -algebra to a measurable space \mathcal{Y} . The probability that W realizes some measurable subset $Z \subseteq \mathcal{Y}$ is given by $P[W \in S] = P[\{w \mid W(w) \in Z\}]$.

The number of false positives is a random variable given by the following theorem..

Theorem 4.1. *Given n negatives, the number of false positives in an approximate set with a false positive rate ε is a random variable denoted by FP_n with a distribution given by*

$$FP_n \sim \text{BIN}(n, \varepsilon) . \quad (4.1)$$

Proof. By axiom 2, the uncertain outcome that a negative element *tests* as positive is a Bernoulli trial with a mean ε . Since there are n such independent and identically distributed trials, the number of false positives is binomially distributed with a mean $n\varepsilon$. \square

The false positive rate ε is an *expectation*. However, the false positive rate of a random approximate set \mathcal{S}^{\pm} parameterized by ε is *uncertain*.

Theorem 4.2. *The false positive rate is the random variable, denoted by \mathcal{E}_n , defined as*

$$\mathcal{E}_n = \frac{FP_n}{n} , \quad (4.2)$$

with an expectation ε , variance $\varepsilon(1 - \varepsilon)/n$, and probability mass function

$$f_{\mathcal{E}_n}(\hat{\varepsilon} \mid \varepsilon) = f_{FP_n}(\hat{\varepsilon}n \mid \varepsilon) . \quad (4.3)$$

over the support $\{\frac{j}{n} \in \mathbb{Q} \mid j \in \{0, \dots, n\}\}$.

Proof. By definition 3.1, the false positive rate is given by the ratio of the number of false positives to the total number of negatives. By theorem 4.1, given that there are n negatives, the number of false positives is a random variable denoted by FP_n . Therefore, the false positive rate, as a function of FP_n , is the random variable $\frac{FP_n}{n}$. The *expected* false positive rate is

$$\mathbb{E}\left[\frac{FP_n}{n}\right] = \frac{1}{n} \mathbb{E}[FP_n] = \varepsilon \quad (a)$$

and its variance is

$$\mathbb{V}\left[\frac{FP_n}{n}\right] = \frac{1}{n^2} \mathbb{V}[FP_n] = \frac{\varepsilon(1 - \varepsilon)}{n} . \quad (b)$$

Finally, $\mathcal{E}_n = FP_n/n$ is a *scaled* transformation of the binomial distribution. Thus, since $FP_n = n\mathcal{E}_n$,

$$f_{\mathcal{E}_n}(\hat{\varepsilon}_n \mid \varepsilon) = f_{FP_n}(n\hat{\varepsilon}_n) . \quad (c)$$

\square

The following corollary immediately follows.

Corollary 4.2.1. *Given n negatives, the number of true negatives in a random approximate set with a false positive rate ε is a random variable denoted by TN_n with a distribution given by*

$$\text{TN}_n = n - \text{FP}_n \sim \text{BIN}(n, 1 - \varepsilon). \quad (4.4)$$

By definition, the true negative rate $\mathcal{N}_n = \text{TN}_n/n = 1 - \varepsilon_n$.

By theorem 4.2, the more negatives there are, the lower the variance.

Corollary 4.2.2. *Given countably infinite negatives, a random approximate set with a false positive rate ε is certain to obtain ε .*

Proof. We know that the *expected* value for each of the random variables in this sequence is ε and the variance is $\varepsilon(1 - \varepsilon)/n$. Immediately, we see that as n increases, the distribution of false positives must become more concentrated around ε . As $n \rightarrow \infty$, the variance goes to 0, i.e., the distribution becomes degenerate with all of the probability mass assigned to the mean. See section A for a more rigorous proof. \square

The fewer negatives, the greater the variance. The maximum possible variance, when $n = 1$ and $\varepsilon = 0.5$, is 0.25, may be used as the most *pessimistic* estimate given a situation where we have no information about the false positive rate ε and the cardinality of the universal set.

A degenerate case is given by letting $n = 0$, corresponding to a random approximate set of the universal set which has no negative elements that can be tested. Respectively, only random *negative* or *positive* approximate sets may be generated for the universal set or empty set.

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

Theorem 4.3. *Given p positives, the uncertain number of false negatives in random approximate sets with a false negative rate ω is modeled as a binomial distributed random variable denoted by FN_p ,*

$$\text{FN}_p \sim \text{BIN}(p, \omega). \quad (4.5)$$

Proof. By axiom 2, the probability that a positive element *tests* as negative is ω . Thus, each test is a Bernoulli trial. Since there are $p = |\mathcal{S}|$ such independent and identically distributed trials with a probability of “success” ω , the number of false negatives is binomially distributed. \square

The false negative rate ω is an *expectation*. However, the false false negative rate of an approximate set \mathcal{S}^\pm parameterized by ω is *uncertain*.

Theorem 4.4. *The false negative rate realizes an uncertain value as given by*

$$\mathcal{W}_p = \frac{\text{FN}_p}{p} \quad (4.6)$$

with a support $\{j/n \mid j = 0, \dots, p\}$, an expectation ω , and a variance $\omega(1 - \omega)/p$.

The proof follows the same logic as the proof for theorem 4.2, except we replace *negatives* with *positives*.

In section 5.1, we consider set-theoretic operations like *complements*. The *complement* operator applied to an approximate set of a set with countably infinite negatives is an approximate set of a set with countably infinite positives.

Corollary 4.4.1. *An approximate set of a set with countably infinite positives has a false negative rate that is certain to obtain ω .*

The proof follows the same logic as the proof for corollary 4.2.2, except we replace *negatives* with *positives*.

The number of true positives is given by the following corollary.

Corollary 4.4.2. *Given p positives, the number of true positives in an approximate set with a false negative rate ω is a random variable denoted by TP_p with a distribution given by*

$$\text{TP}_p \sim \text{BIN}(p, \tau). \quad (4.7)$$

By definition, the true positive rate is given by $\mathcal{T}_p = 1 - \mathcal{W}_p$.

The proof follows the same logic as the proof for theorem 4.2.

Many other properties of random approximate sets follow from these distributions. For instance,

$$|\mathcal{A}_\varepsilon^+| = \text{TP}_p + \text{FP}_n, \quad (4.8)$$

which has an expectation of $n\varepsilon + p\tau$ and variance of $n\varepsilon(1 - \varepsilon) + p\tau(1 - \tau)$, which is the generalization of the binomial distribution known as the *Poisson binomial distribution*.

If we do not know p , the cardinality \mathcal{A} , but have observed $\mathcal{A}^\pm = \mathcal{B}$, then \mathcal{B} has a cardinality that tends to be centered around $u\varepsilon + p(\tau - \varepsilon)$ where u is the cardinality of the universal set. Solving for p yields a method of moments estimator

$$\hat{p} = \frac{|\mathcal{B}| - u\varepsilon}{\tau - \varepsilon}. \quad (4.9)$$

If the universal set \mathcal{U} is infinite, then this estimator is undefined.

4.1 Asymptotic limits

The false positive and false negative rates are a function of the cardinality of the objective and universal sets. The limiting distributions for the false positive and true positive rates are given by the following theorems.

Theorem 4.5. *By theorem 4.2, the uncertain false positive rate \mathcal{E}_n converges in distribution to the normal distribution with a mean ε and a variance $\varepsilon(1 - \varepsilon)/n$, written*

$$\mathcal{E}_n \xrightarrow{d} \mathcal{N}(\varepsilon, \varepsilon(1 - \varepsilon)/n). \quad (4.10)$$

Similarly, by ??, the uncertain true positive rate of an approximate set of p positives, denoted by \mathcal{T}_p , converges in distribution to the normal distribution with a mean τ and a variance $\tau(1 - \tau)/p$, written

$$\mathcal{T}_p \xrightarrow{d} \mathcal{N}(\tau, \tau(1 - \tau)/p). \quad (4.11)$$

Proof. By ?? in the proof of corollary 4.2.2, given n negatives, the false positive rate is

$$\mathcal{E}_n = \frac{X_1}{n} + \dots + \frac{X_n}{n}, \quad (a)$$

where X_1, \dots, X_n are n independent Bernoulli trials each with a mean ε and a variance $\varepsilon(1 - \varepsilon)$. Therefore, by the central limit theorem, \mathcal{E}_n converges in distribution to a normal distribution with a mean ε and a variance $\varepsilon(1 - \varepsilon)/n$. The proof for the true positive rate follows the same logic. \square

By eqs. (4.10) and (4.11),

$$\mathcal{N}_n \xrightarrow{d} \mathcal{N}(1 - \varepsilon, \varepsilon(1 - \varepsilon)/n) \text{ and } \mathcal{W}_n \xrightarrow{d} \mathcal{N}(1 - \tau, \tau(1 - \tau)/p). \quad (4.12)$$

The random approximate set model is the *maximum entropy* probability distribution for the indicated false positive and true positive rates, e.g., any estimated α -confidence intervals are the largest intervals possible for the indicated α and therefore represent a worst-case uncertainty.

If we generate an approximate set, the uncertain false positive and true positive rates realize certain values, i.e., $\mathcal{E}_n = \hat{\varepsilon}$ and $\mathcal{T}_p = \hat{\tau}$. If the sample space is countably infinite, the distribution is degenerate, e.g., $\mathcal{E}_n = \varepsilon$ with probability 1. However, for finite sample spaces, the outcomes are uncertain. If these outcomes can be *observed*, e.g., it is not too costly to compute, the exact values $\hat{\varepsilon}$ and $\hat{\tau}$ may be recorded. If these outcomes cannot be observed, e.g., it is too costly to compute or the information to compute $\hat{\varepsilon}$ or $\hat{\tau}$ is not available, we may use the probabilistic model to inform us about the distribution of false positive rates.

Confidence intervals that contain the true false positive rate $\hat{\varepsilon}$ and the true true positive rate $\hat{\tau}$ are given by the following corollaries.

Theorem 4.6. *Given a random approximate set parameterized by ε and τ , asymptotic $\alpha \cdot 100\%$ confidence intervals for the false positive rate and true positive rate are respectively*

$$\varepsilon \pm \sqrt{\frac{\varepsilon(1 - \varepsilon)}{n}} \Phi^{-1}(\alpha/2) \quad (4.13)$$

and

$$\tau \pm \sqrt{\frac{\tau(1 - \tau)}{p}} \Phi^{-1}(\alpha/2), \quad (4.14)$$

where $\Phi^{-1}: [0, 1] \mapsto \mathbb{R}$ is the inverse cumulative distribution function of the standard normal.

As a worst-case (maximum uncertainty), we may let $n = p = 1$ in eqs. (4.13) and (4.14).

5 Probability distributions of parameters for *higher-order* approximations

The false negative rate is a random variable that is a mixture? of binomially distributed random variables and likewise for the false positive rate. We characterize their error rates by their expected values and variance and show their normal approximation when the sampling distribution of these rates are important.

5.1 Compositions of random approximate sets

The union of two first-order random approximate sets is given by the following theorem.

Theorem 5.1. *The union of first-order random approximate sets \mathcal{A}^{ϵ_1} and \mathcal{B}^{ϵ_2} is up to a third-order random approximate set with a random error rate given by*

$$\Delta = \sum_{i=1}^3 \alpha_i X_i \quad (5.1)$$

with an expectation

$$\epsilon = \alpha_1(1 - \epsilon_1)\epsilon_2 + \alpha_2(1 - \epsilon_2)\epsilon_1 + (1 - \alpha_1 - \alpha_2 - \alpha_3)\epsilon_1\epsilon_2, \quad (5.2)$$

a false negative rate

$$\omega = \epsilon_1\epsilon_2 \quad (5.3)$$

and a false positive rate

$$\varepsilon = \epsilon_1\epsilon_2, \quad (5.4)$$

where $\alpha_1 := \frac{|A \setminus B|}{|\mathcal{U}|}$, $\alpha_2 := \frac{|B \setminus A|}{|\mathcal{U}|}$, and $\alpha_3 := \frac{|A \cap B|}{|\mathcal{U}|}$.

The set of n -arity operations on set \mathcal{X} is given by $\mathcal{X}^n \mapsto \mathcal{X}$. In this section, we consider binary operations on $2^{\mathcal{U}}$, like $\cup: 2^{\mathcal{U}} \times 2^{\mathcal{U}} \mapsto 2^{\mathcal{U}}$, and the result of providing *random approximate sets* as input.

Given the Boolean algebra $(2^{\mathcal{U}}, \cap, \cup, \bar{\cdot}, \emptyset, \mathcal{U})$ we derive the random approximate sets that result from the union (join) or complement of random approximate sets. Since unions and complements form a *complete basis* for Boolean algebras, we may express other Boolean operations as a composition of these two operations, e.g., $\mathcal{A}^{\pm} \setminus \mathcal{B}^{\pm} = \overline{\mathcal{A}^{\pm} \cup \mathcal{B}^{\pm}}$.

The random approximate sets that result from union operations on random approximate sets are given by the following theorems.

Theorem 5.2. *The union of two random approximate sets respectively with true negative rates η_1 and η_2 is a random approximate set with a true negative rate $\eta_1\eta_2$.*

Proof. Suppose we have two sets \mathcal{A} and \mathcal{B} with false positive rates ε_1 and ε_2 . The false positive rate ε of $\mathcal{A}^{\pm} \cup \mathcal{B}^{\pm}$ is a probability conditioned on a negative for $\mathcal{A} \cup \mathcal{B}$ being a positive for $\mathcal{A}^{\pm} \cup \mathcal{B}^{\pm}$.

Switching to the Boolean vector representation, suppose we randomly select an element from the universe, denoted by x_j , such that $\neg A_j \vee \neg B_j$ is true.

The expected false positive rate of the union is defined by the probability

$$\varepsilon = \mathbb{P}[\mathcal{A}^{\pm} \cup \mathcal{B}^{\pm} \mid B_1 \cap B_2]. \quad (a)$$

By DeMorgan's law, the union of sets is the complement of the intersection of their complements. That is,

$$A_1 \cup A_2 \equiv (A'_1 \cap A'_2)' \quad (b)$$

and thus

$$\varepsilon = \mathbb{P}[(A'_1 \cap A'_2)' \mid B_1 \cap B_2]. \quad (c)$$

Since either an event or the *complement* of the event is certain to occur, $\mathbb{P}[E] + \mathbb{P}[E'] = 1$, the above equation may be rewritten as

$$\varepsilon = 1 - \mathbb{P}[A'_1 \cap A'_2 \mid B_1 \cap B_2]. \quad (d)$$

Since A'_1 and A'_2 are independent,

$$\varepsilon = 1 - \mathbb{P}[A'_1 \mid B_1 \cap B_2] \mathbb{P}[A'_2 \mid B_1 \cap B_2]. \quad (e)$$

Since A_1 is conditionally independent of B_2 and A_2 is conditionally independent of B_1 , we may rewrite the above equation as

$$\varepsilon = 1 - \mathbb{P}[A'_1 \mid B_1] \mathbb{P}[A'_2 \mid B_2]. \quad (f)$$

A_j denotes $X \in \mathcal{S}^j$, therefore A'_j denotes $X \notin \mathcal{S}^j$. Substituting the definition of A'_1 , A'_2 , B_1 , and B_2 into the above equation gives

$$\varepsilon = 1 - \mathbb{P}[X \in \mathcal{A}^{\pm} \mid X \notin \mathcal{A}] \mathbb{P}[X \notin \mathcal{B}^{\pm} \mid X \notin \mathcal{B}]. \quad (g)$$

By definition, $P[X \notin \mathcal{A}^\pm \mid X \notin \mathcal{A}]$ is the true negative rate η_1 and likewise for \mathcal{B}^\pm . Thus,

$$\varepsilon = 1 - \eta_1 \eta_2. \quad (\text{h})$$

□

The limiting probability distribution of the uncertain true negative rate of the union of \mathcal{A}^\pm and \mathcal{B}^\pm is thus

$$\mathcal{N}_n \sim \mathcal{N}\left(\eta_1 \eta_2, \frac{\eta_1 \eta_2 (1 - \eta_1 \eta_2)}{n}\right) \quad (5.5)$$

where the number of negatives $n = u - |\mathcal{A} \cup \mathcal{B}| \leq u$, which is a value between 0 and u . Since this is a limiting distribution, presumably n is large, and as $n \rightarrow \infty$ the distribution converges in probability to $\tau_1 \tau_2$.

Generally, the number of negatives n or positives p is not known, and so this serves a more analytic function, i.e., given around n negatives, what true negative rate η provides the desired level of confidence that the true negative rate will not realize a value less than some specified value?

Theorem 5.3. *The union of $\mathcal{A}_{\eta_1}^{\omega_1}$ and $\mathcal{B}_{\eta_2}^{\omega_2}$ is a random approximate set with an expected false negative rate*

$$\omega = \alpha_1 \omega_1 \eta_2 + \alpha_2 \eta_1 \omega_2 + (1 - \alpha_1 - \alpha_2) \omega_1 \omega_2, \quad (5.6)$$

where

$$\begin{aligned} 0 \leq \alpha_1 &= \frac{|\mathcal{A} \setminus \mathcal{B}|}{|\mathcal{A} \cup \mathcal{B}|}, \\ 0 \leq \alpha_2 &= \frac{|\mathcal{B} \setminus \mathcal{A}|}{|\mathcal{A} \cup \mathcal{B}|}, \end{aligned} \quad (5.7)$$

$$\alpha_1 + \alpha_2 \leq 1.$$

See section C for a proof of theorem 5.3.

The complement of an approximate set is given by the following theorem.

Theorem 5.4. *The complement of a random approximate set with a false positive rate ε and false negative rate ω is an approximate set with a false positive rate ω and a false negative rate ε .*

Proof. The false positives in an approximate set are false negatives in its complement; likewise, the false negatives in an approximate set are the false positives in its complement set. □

Remark. *Consider a sequence $\mathcal{A}_{\varepsilon_1}^+, \dots, \mathcal{A}_{\varepsilon_n}^+$. Any subsequence contains strictly less information about \mathcal{A} . That is, positive approximate sets are strictly additive, and as $n \rightarrow \infty$, $\cap_{i=1}^n \mathcal{A}_{\varepsilon_i}^+$ converges almost surely to \mathcal{A} .³ △*

6 Distribution of binary classification measures

Suppose we have some other function $g: 2^{\mathcal{X}} \mapsto \mathcal{Y}$ that is not a *constant*, then the composition $g \circ \text{id}_\varepsilon^r$ is some probability distribution over the codomain \mathcal{Y} . That is, $(g \circ \text{id}_\varepsilon^r)(\mathcal{A})$ is a *random variable*.

Example 2 *Let $f: 2^{\{0,1\}} \mapsto \{0,1\}$ be defined as*

$$f(\mathcal{A}) := \begin{cases} 1 & \mathcal{A} \in \{\{1\}, \{0,1\}\}, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{a})$$

The composition $f \circ \text{id}_{.75}^{.25}$ generates Bernoulli distribution random variables, e.g., $(f \circ \text{id}_{.75}^{.25})(\{0\}) \sim \text{BER}(0.25)$.

We consider several classes of functions and the distributions induced by replacing the inputs with random approximate sets, e.g., operators like set-union or binary performance measures like positive predictive value. In section 9, we consider a more sophisticated example in Boolean search where queries map to *random approximate result sets*.

In the approximate set model, the distribution of random variables like the false positive, false negatives, true positives, and true negative rates are given respectively by parameters ε , ω , τ , and η . These parameters belong to a more general class of *binary performance measures*.

The above parameters are statements about the distribution of random approximate sets given corresponding objective sets of interest, e.g.,

$$P[\mathbb{1}_{\mathcal{A}_\varepsilon^r}(x) \mid \mathbb{1}_{\mathcal{A}}(x)] = \tau. \quad (6.1)$$

The accuracy of *predictions* about objective sets given a corresponding approximate set is usually the more relevant performance measure. The *positive predictive value* is given by the following definition.

³Likewise, for negative approximate sets, as $n \rightarrow \infty$, $\cup_{i=1}^n \mathcal{A}_{\omega_i}^-$ converges almost surely to \mathcal{A} .

Definition 6.1. The positive predictive value is a performance measure defined as

$$\text{ppv} = \frac{t_p}{t_p + f_p} \quad (6.2)$$

where t_p is the number of true positives and f_p is the number of false positives.

The positive predictive value of random approximate sets is a random variable given by the following theorem.

Theorem 6.1. Given n negatives, p positives, and a random approximate set with false positive and true positive rates ε and τ respectively, the positive predictive value is a random variable

$$\text{PPV} = \frac{\text{TP}_p}{\text{TP}_p + \text{FP}_n} \quad (6.3)$$

with an expectation given approximately by

$$\text{ppv}(\tau, \varepsilon, p, n) \approx \frac{\bar{t}_p}{\bar{t}_p + \bar{f}_p} + \frac{\bar{t}_p \sigma_{f_p}^2 - \bar{f}_p \sigma_{t_p}^2}{(\bar{t}_p + \bar{f}_p)^3}, \quad (6.4)$$

where $\bar{t}_p = p\tau$ is the expected true positive frequency, $\bar{f}_p = n\varepsilon$ is the expected false positive frequency, $\sigma_{t_p}^2 = (1 - \tau)\bar{t}_p$ is the variance of the true positive frequency, and $\sigma_{f_p}^2 = (1 - \varepsilon)\bar{f}_p$ is the variance of false positive frequency.

See section B for a proof of theorem 6.1.

We make the following observations about eq. (6.4):

1. For sufficiently large approximate sets, $\text{ppv} \approx \bar{t}_p / (\bar{t}_p + \bar{f}_p)$.
2. If $\varepsilon \neq 0$, as $n \rightarrow \infty$, $\text{ppv} \rightarrow 0$.
3. As $\varepsilon \rightarrow 0$, $\text{ppv} \rightarrow 1$.

Accuracy is given by the following definition.

Definition 6.2. The accuracy is the proportion of true results (both true positives and true negatives) in the universe of positives and negatives, $(t_p + t_n) / (p + n)$, where t_p , t_n , p , and n are respectively the number of true positives, true negatives, positives, and negatives.

The expected accuracy is given by the following theorem.

Theorem 6.2. Given p positives and n negatives, a random approximate set with an expected false positive rate ε and an expected true positive rate τ is a random variable given by

$$\text{ACC}_{p+n} = \lambda \mathcal{T}_p + (1 - \lambda) \mathcal{N}_n. \quad (6.5)$$

has an expected accuracy

$$\text{acc}(\tau, \varepsilon, n, p) = \lambda \tau + (1 - \lambda) \eta \quad (6.6)$$

with a variance

$$\frac{\lambda \omega \tau + (1 - \lambda) \varepsilon \eta}{p + n}, \quad (6.7)$$

where $\lambda = p / (p + n)$.

Proof. Suppose there the u elements in the universe can be partitioned into p positives and n negatives. An approximate set \mathcal{S}^\pm with a false positive rate ε and false negative rate ω has an uncertain accuracy

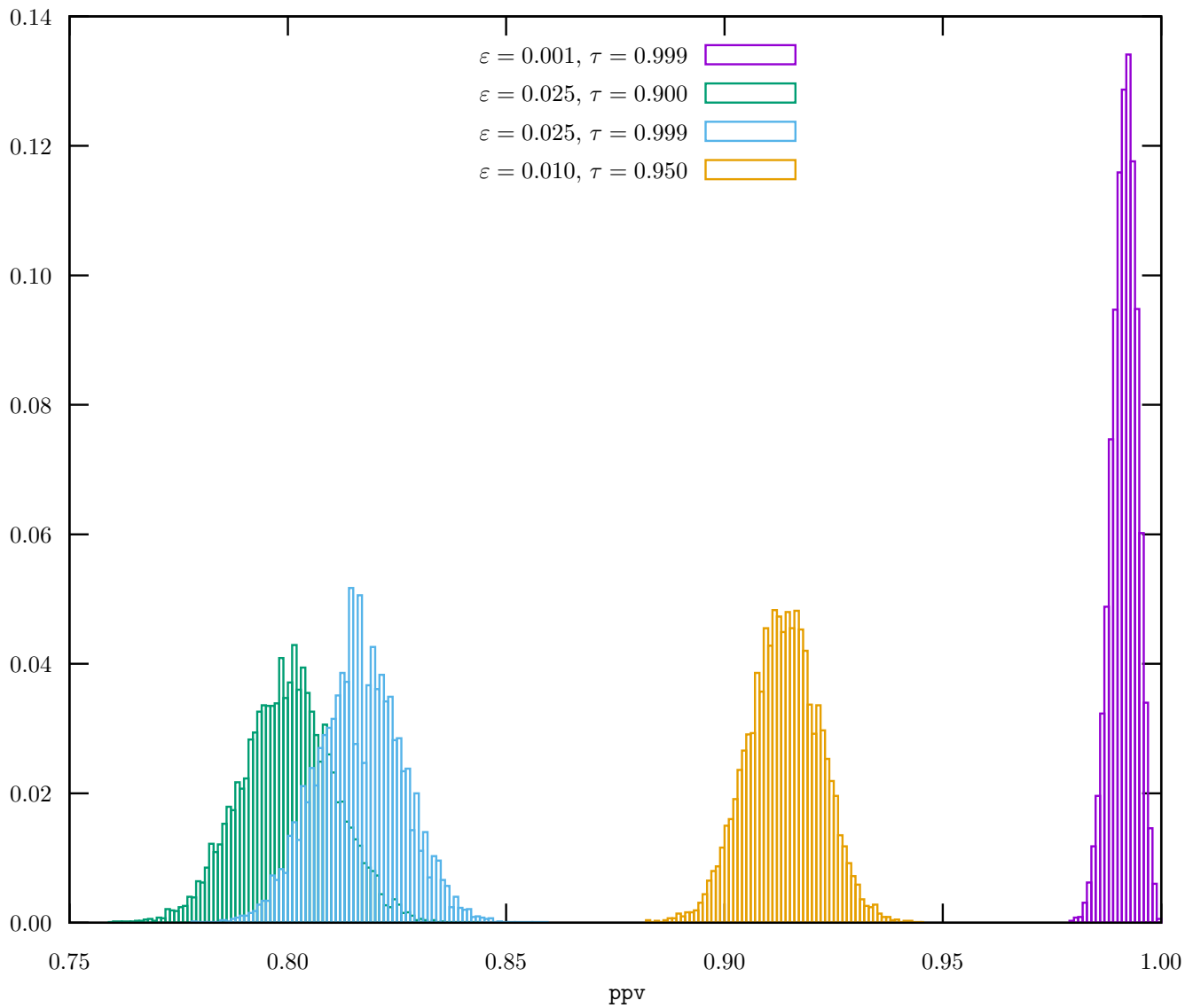
$$\text{ACC}_{p+n} = \frac{\text{TP}_p + \text{TN}_n}{p + n}. \quad (\text{a})$$

The expected accuracy is given by the expectation

$$\text{E}[\text{ACC}_{p+n}] = \text{E}\left[\frac{\text{TP}_p + \text{TN}_n}{p + n}\right] \quad (\text{b})$$

$$= \frac{p(1 - \omega) + n(1 - \varepsilon)}{p + n}. \quad (\text{c})$$

Figure 2: Relative frequency of positive predictive values for several different parameterizations of the false positive and true positive rates given $n = 900$ negatives and $p = 100$ positives.



Noting that $n/(p+n) = 1 - p/(p+n)$ and letting $\lambda = p/(p+n)$,

$$E[\text{ACC}_{p+n}] = \lambda(1 - \omega) + (1 - \lambda)(1 - \varepsilon). \quad (\text{d})$$

The variance

$$V[\text{ACC}_{p+n}] = V\left[\frac{\text{TP}_p}{p+n}\right] + V\left[\frac{\text{TN}_n}{p+n}\right] \quad (\text{e})$$

$$= \frac{1}{(p+n)^2} V[\text{TP}_p] + \frac{1}{(p+n)^2} V[\text{TN}_n] \quad (\text{f})$$

$$= \frac{p\omega(1-\omega)}{(p+n)^2} + \frac{n\varepsilon(1-\varepsilon)}{(p+n)^2} \quad (\text{g})$$

$$= \frac{\lambda\omega\tau + (1-\lambda)\varepsilon\eta}{p+n}. \quad (\text{h})$$

□

Negative predictive value is given by the following definition.

Definition 6.3.

$$\text{npv} = \frac{t_n}{t_n + f_n} \quad (\text{6.8})$$

where t_n and f_n are respectively the number of true negatives and false negatives

The expected negative predictive value is given by the following theorem.

Theorem 6.3. *Given p positives, n negatives, and a random approximate set with false positive and true positive rates ε and τ respectively, the negative predictive value is a random variable*

$$\text{NPV} = \frac{\text{TN}_n}{\text{TN}_n + \text{FN}_p} \quad (\text{6.9})$$

with an expectation given approximately by

$$\text{npv}(\tau, \varepsilon, p, n) \approx \frac{\bar{t}_n}{\bar{t}_n + \bar{f}_n} + \frac{\bar{t}_n\sigma_{f_n}^2 - \bar{f}_n\sigma_{t_n}^2}{(\bar{t}_n + \bar{f}_n)^3}, \quad (\text{6.10})$$

where $\bar{t}_n = n(1 - \varepsilon)$ is the expected true negative frequency, $\bar{f}_n = p(1 - \tau)$ is the expected false negative frequency, $\sigma_{t_n}^2 = \varepsilon\bar{t}_n$ is the variance of the true negative frequency, and $\sigma_{f_n}^2 = \tau\bar{f}_n$ is the variance of the false negative frequency.

The proof for theorem 6.3 follows the same pattern as the proof for theorem 6.1.

Youden's J statistic is a measure of the performance of a binary test, defined as

$$J = \frac{t_p}{t_p + f_n} + \frac{t_n}{t_n + f_p} - 1, \quad (\text{6.11})$$

with a range $[0, 1]$. In the case of the random approximate set model, J is a random variable

$$J = \mathcal{T}_p - \mathcal{E}_n, \quad (\text{6.12})$$

which has an expectation

$$E[J] = \tau - \varepsilon. \quad (\text{6.13})$$

Table 2 may be used to reparameterize an approximate set.

Example 3 *Suppose we seek a positive approximate set with an expected accuracy γ . By table 2,*

$$\gamma = \text{acc}(\varepsilon, \omega = 0, \lambda) = 1 - \varepsilon(1 - \lambda). \quad (\text{a})$$

Solving for ε in terms of γ yields the result

$$\varepsilon(\gamma, \lambda) = \frac{1 - \gamma}{1 - \lambda} \quad (\text{b})$$

subject to $0 \leq \lambda \leq \gamma \leq 1$ and $\lambda < 1$. Under this parameterization of the positive approximate set, λ must be known (or estimated). Note that if $\lambda = 1$ then $\varepsilon(\gamma, \lambda = 1)$ is undefined as expected, but as λ goes to 1, $\varepsilon(\cdot; \lambda)$ goes to 1 and γ goes to 1, which logically follows since if there are no negatives, there can be no false positives.

Table 2: Various *expected* performance measures.

measure	parameter	expected value
true positive rate	$\text{tpr}(\tau)$	τ
false positive rate	$\text{fpr}(\varepsilon)$	ε
false negative rate	$\text{fnr}(\tau)$	$1 - \tau$
true negative rate	$\text{tnr}(\varepsilon)$	$1 - \varepsilon$
accuracy	acc	eq. (6.6)
positive predictive value	ppv	eq. (6.4)
negative predictive value	npv	eq. (6.10)
false discovery rate	fdr	$1 - \text{ppv}$
false omission rate	for	$1 - \text{npv}$

7 Uncertain rate distortions

We may not be certain about the *expected* false positive and true positive rates, i.e., we may only have the joint distribution of \mathcal{A}^\pm , \mathcal{E} , and \mathcal{W} .

7.1 First-order model

The easiest case to analyze is the *first-order* random approximate set model. Suppose we are interested in the distribution of \mathcal{A}^\pm given \mathbf{R} has p positives and n negatives. Since we are primarily interested in the distribution of false positives and true positives (or their corresponding rates), we consider the related random tuple $\langle \text{FP}_n, \text{TP}_p, \mathcal{T}, \mathcal{E} \rangle$ which, assuming \mathcal{E} and \mathcal{T} are independent, has a joint probability density function given by

$$f_{\text{TP}_p, \text{FP}_n, \mathcal{T}, \mathcal{E}}(t, f, \tau, \varepsilon) = f_{\text{TP}_p, \mathcal{T}}(t, \tau) f_{\text{FP}_n, \mathcal{E}}(f, \varepsilon) \quad (7.1)$$

where

$$f_{\text{TP}_p, \mathcal{T}}(t, \tau) = f_{\text{TP}_p | \mathcal{T}}(t | \tau) f_{\mathcal{T}}(\tau), \quad (7.2)$$

$$f_{\text{FP}_n, \mathcal{E}}(f, \varepsilon) = f_{\text{FP}_p | \mathcal{E}}(f | \varepsilon) f_{\mathcal{E}}(\varepsilon). \quad (7.3)$$

When we *marginalize* over the true positives, we get the result

$$\begin{aligned} f_{\text{TP}_p}(t) &= \int_0^1 f_{\text{TP}_p | \mathcal{T}}(t | \tau) f_{\mathcal{T}}(\tau) d\tau \\ &= \int_0^1 \binom{p}{t} \tau^t (1 - \tau)^{p-t} f_{\mathcal{T}}(\tau) d\tau. \end{aligned} \quad (7.4)$$

If all the probability mass for \mathcal{T} is assigned to a particular point τ , the probability mass function simplifies to

$$f_{\text{TP}_p}(t) = \binom{p}{t} \tau^t (1 - \tau)^{p-t}, \quad (7.5)$$

which is probability mass function of a binomial distribution.

The simplest kind of uncertainty is given by a disjoint set of intervals, in which the true expected rate is uniformly distributed across the support.

Definition 7.1. *An interval is a convex set of real numbers. We denote by $[x] = [\underline{x}, \bar{x}]$ an interval with a lower-bound \underline{x} and an upper-bound \bar{x} .*

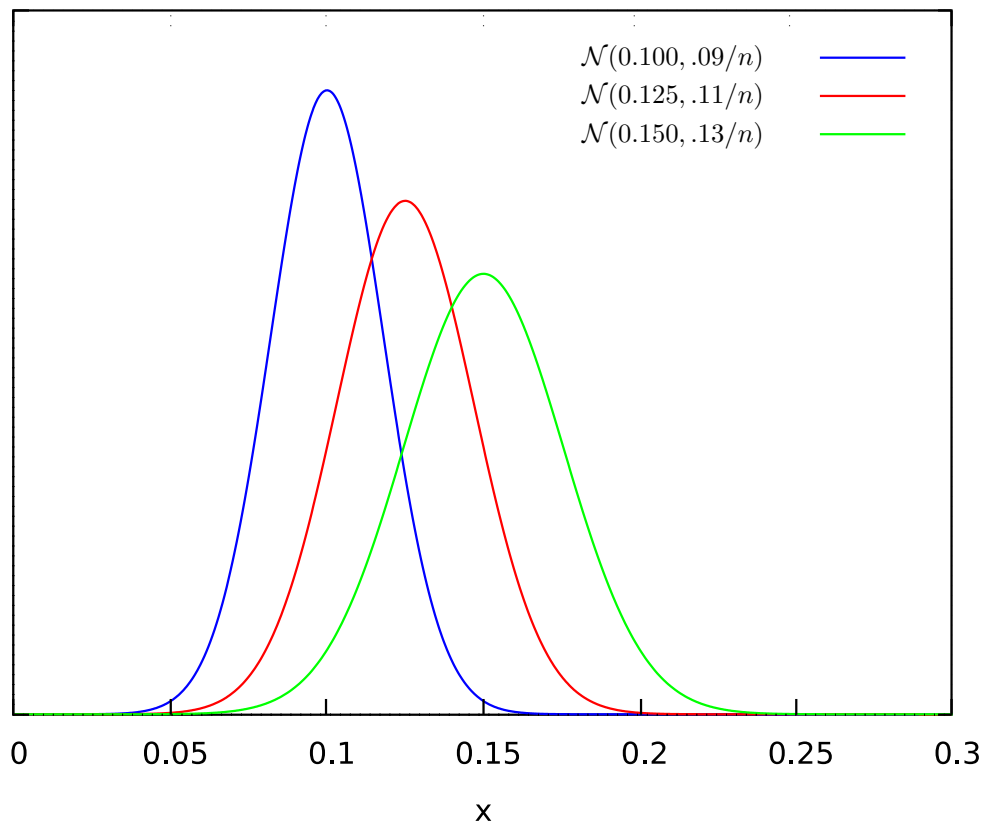
A further simplification comes from mapping the disjoint set of intervals to the smallest interval that *spans* all of them.

Definition 7.2. *Given a disjoint interval set \mathcal{X} , $\text{span}(\mathcal{X})$ maps to an interval with lower and upper bounds that are the lower and upper bounds of \mathcal{X} .*

A confidence interval, for instance, may be specified in this notation. Here, however, we consider an algebra for interval arithmetic and put it to use quantifying our ignorance about the distribution of parameters after, for instance, a union operation.

The performance measures summarized by table 2 depend upon the false positive rate ε , false negative rate ω , and proportion of positives λ being *known*. Any parameters that are not known with certainty may be replaced in the above table

Figure 3: Relative frequency of positive predictive values for several different parameterizations of the false positive and true positive rates given $n = 900$ negatives and $p = 100$ positives.



by intervals that (are assumed to) contain the expected value. As a consequence, the performance measure will also be an interval.

Maximum uncertainty is when the parameter value is in the interval $[0, 1]$, e.g., $[\lambda] = [0, 1]$, and *minimum* uncertainty is when the parameter is some value in the degenerate interval $[x, x]$, e.g., $[\varepsilon] = [.2, .2]$. The more certain—the smaller the width of the intervals—the more certain the performance measure.

When using interval arithmetic, the *dependency problem* can lead to overly pessimistic bounds. In our case, the formulae are simple enough to ensure dependencies are satisfied. We show the results of an uncertain proportion of positives $[\lambda]$ for the *accuracy* measure in the following example.

Example 4 *Suppose we wish to determine the expected accuracy given that the proportion of positives is known to be some value in the interval $[\lambda]$. Then, the expected accuracy is some value in the interval*

$$\text{acc}([\varepsilon], [\omega]; [\lambda]) = \left[\begin{aligned} &f(\bar{\varepsilon}, \bar{\omega})(1 - \bar{\omega}) + (1 - f(\bar{\varepsilon}, \bar{\omega}))(1 - \bar{\varepsilon}), \\ &f(\underline{\omega}, \underline{\varepsilon})(1 - \underline{\omega}) + (1 - f(\underline{\omega}, \underline{\varepsilon}))(1 - \underline{\varepsilon}) \end{aligned} \right], \quad (\text{a})$$

where $f(x, y) = \bar{\lambda}[x < y] + \underline{\lambda}[y \leq x]$. If we have complete ignorance about λ then $[\lambda] = [0, 1]$. As a special case, if we have complete ignorance about lambda and $\omega = 0$ (positive approximate set), then $\text{acc}([\varepsilon], 0; [0, 1]) = [1 - \bar{\varepsilon}, 1]$.

However, the expected rates may not be known, e.g., the values of α_1 and α_2 in eq. (5.6) may not be known. Alternatively, we may not be interested in the *expected* value, but the smallest set of values such that with probability $1 - \alpha$ the true rate realizes some value in the set, which is typically an *interval*, i.e., a confidence interval.

Intervals represent an uncertainty and they manifest themselves in two independent ways. The common notion of the *confidence interval* is a product of the probabilistic model, i.e., the realized true positive rate $\hat{\tau}$, which is normally centered around the expected true positive rate τ as discussed in section 4.1. We may use *interval arithmetic* and replace point values interval values, point values being a degenerate case. Basic interval arithmetic is presented in [?].

A set sampled from $\mathcal{A}^\pm(\varepsilon, \tau)$ is an approximate set such that the $(1 - \alpha)\%$ asymptotic confidence interval for the false negative and false positive rates given respectively by

$$[\omega] =? \quad (7.6)$$

and

$$[\varepsilon] =?. \quad (7.7)$$

By ??, $\mathcal{A}^\pm \cup \mathcal{B}^\pm$, the observation $\mathcal{A}^\pm(\varepsilon, \tau) = \mathcal{A}$ is an approximate set with a false negative rate

$$\hat{\omega} \in [\omega_1](1 - [\varepsilon_2]) \cup [\omega_2](1 - [\varepsilon_1]) \cup [\omega_1][\omega_2] \quad (7.8)$$

and a false positive rate

$$\hat{\varepsilon} \in 1 - (1 - [\varepsilon_1])(1 - [\omega_2]). \quad (7.9)$$

Equation (7.10) represents a disjoint set of intervals. However, we are only interested in the best and worst case of the false negative rate. Thus, we map the disjoint set to a *minimum width* interval that contains every point in the disjoint set.

Definition 7.3. *Given a set \mathcal{X} , $\text{span}(\mathcal{X})$ maps to an interval with lower and upper bounds that are the lower and upper bounds of \mathcal{X} .*

Theorem 7.1. *The union of two approximate sets with uncertain false negative rates $[\omega_1]$ and $[\omega_2]$ and uncertain false positive rates $[\varepsilon_1]$ and $[\varepsilon_2]$ is an approximate set with an uncertain false negative rate*

$$\begin{aligned} [\omega] &= \text{span}([\omega_1](1 - [\varepsilon_2]) \cup [\omega_2](1 - [\varepsilon_1]) \cup [\omega_1][\omega_2]) \\ &= \left[\begin{aligned} &\min\{\underline{\omega}_1(1 - \bar{\varepsilon}_2), \underline{\omega}_2(1 - \bar{\varepsilon}_1), \underline{\omega}_1\underline{\omega}_2\}, \\ &\max\{\bar{\omega}_1(1 - \underline{\varepsilon}_2), \bar{\omega}_2(1 - \underline{\varepsilon}_1), \bar{\omega}_1\bar{\omega}_2\} \end{aligned} \right] \end{aligned} \quad (7.10)$$

and an uncertain true negative rate

$$\begin{aligned} [\eta] &= [\eta_1][\eta_2] \\ &= [\eta_1\underline{\eta}_2, \bar{\eta}_1\bar{\eta}_2] \\ &= 1 - (1 - [\varepsilon_1])(1 - [\varepsilon_2]) \\ &= [\underline{\varepsilon}_1 + \underline{\varepsilon}_2 - \underline{\varepsilon}_1\underline{\varepsilon}_2, \bar{\varepsilon}_1 + \bar{\varepsilon}_2 - \bar{\varepsilon}_1\bar{\varepsilon}_2]. \end{aligned} \quad (7.11)$$

Proof. By ??, the false positive rate of $\mathcal{A}^\pm \cup \mathcal{B}^\pm$ is

$$[\varepsilon] = [\varepsilon_1] + [\varepsilon_2] - [\varepsilon_1][\varepsilon_2]. \quad (\text{a})$$

and the false negative rate is

$$\begin{aligned} [\hat{\omega}] &= \alpha_1 [\hat{\omega}_1] (1 - [\varepsilon_2]) + \alpha_2 [\hat{\omega}_2] (1 - [\varepsilon_1]) \\ &\quad + (1 - \alpha_1 - \alpha_2) (1 - [\varepsilon_1] + [\hat{\omega}_2][\varepsilon_1]), \end{aligned} \quad (\text{5.6 revisited})$$

where $\alpha_1, \alpha_2 \geq 0$ and $\alpha_1 + \alpha_2 \leq 1$. Thus, to maximize (minimize) this equation, we simply need to put all of the *weight* into the largest (smallest) term. \square

Theorem 7.2. *The complement of an approximate set with a false negative rate $[\omega]$ and false positive rate $[\varepsilon]$ is an approximate set with a false negative rate $[\varepsilon]$ and false positive rate $[\omega]$.*

Since any set-theoretic composition is reducible to a combination of unions and complements, we may use theorems 7.1 and 7.2 to compute the bounds for any set-theoretic composition of approximate sets. See ?? for a summary of a several well-known operations.

Table 3: The smallest intervals that contain the false positive and false negative rates of the approximate sets that result from the corresponding set-theoretic operations on approximate sets \mathcal{A}^\pm ($[\omega_1], [\varepsilon_1]$) and \mathcal{B}^\pm ($[\omega_2], [\varepsilon_2]$).

op	param	interval
$\mathcal{A}^\pm \cup \mathcal{B}^\pm$	$[\varepsilon]$	$1 - (1 - [\varepsilon_1])(1 - [\varepsilon_2])$
	$[\omega]$	$\text{span}([\omega_1](1 - [\varepsilon_2]) \cup [\omega_2](1 - [\varepsilon_1]) \cup [\omega_1][\omega_2])$
$\mathcal{A}^\pm \cap \mathcal{B}^\pm$	$[\varepsilon]$	$\text{span}([\omega_1](1 - [\omega_2]) \cup [\varepsilon_2](1 - [\omega_1]) \cup [\omega_1][\varepsilon_2])$
	$[\omega]$	$1 - (1 - [\omega_1])(1 - [\omega_2])$
$\mathcal{A}^\pm \setminus \mathcal{B}^\pm$	$[\varepsilon]$	$\text{span}([\omega_1](1 - [\varepsilon_2]) \cup [\omega_2](1 - [\omega_1]) \cup [\omega_1][\omega_2])$
	$[\omega]$	$[\underline{\omega}_1 + \underline{\varepsilon}_2(1 - \bar{\omega}_1), \bar{\omega}_1 + \bar{\varepsilon}_2(1 - \underline{\omega}_1)]$
$\overline{\mathcal{A}^\pm}$	$[\varepsilon]$	$[\omega_1]$
	$[\omega]$	$[\varepsilon_1]$

Example 5 *Suppose we have three sets \mathcal{A} , \mathcal{B} , and \mathcal{C} and consider the random approximate set*

$$\mathcal{D}_\varepsilon^r = (\mathcal{A}_\varepsilon^+ \cap \mathcal{B}_\varepsilon^+) \setminus \mathcal{C}_\varepsilon^+. \quad (\text{a})$$

The intersection of \mathcal{A}^+ and \mathcal{B}^+ is an approximate set

$$\mathcal{A}^+ \cap \mathcal{B}_{[\varepsilon]}^+ = \overline{\overline{\mathcal{A}^\pm} \cup \overline{\mathcal{B}^\pm}}. \quad (\text{b})$$

Table 4: The tightest intervals that contain the false positive and false negative rates of the positive or negative approximate sets that result from the corresponding set-theoretic operations.

(a) $\mathcal{A}^+([\varepsilon_1])$ and $\mathcal{A}^+([\varepsilon_2])$.			(b) $\mathcal{A}^-([\omega_1])$ and $\mathcal{B}^-([\omega_2])$.		
op	param	interval	op	param	interval
$\mathcal{A}^+ \cup \mathcal{B}^+$	$[\varepsilon]$	$1 - (1 - [\varepsilon_1])(1 - [\varepsilon_2])$	$\mathcal{A}^- \cup \mathcal{B}^-$	$[\omega]$	$[\underline{\omega}_1 \underline{\omega}_2, \max(\bar{\omega}_1, \bar{\omega}_2)]$
$\mathcal{A}^+ \cap \mathcal{B}^+$	$[\varepsilon]$	$[\underline{\varepsilon}_1 \underline{\varepsilon}_2, \max(\bar{\varepsilon}_1, \bar{\varepsilon}_2)]$	$\mathcal{A}^- \cap \mathcal{B}^-$	$[\omega]$	$1 - (1 - [\omega_1])(1 - [\omega_2])$
$\mathcal{A}^+ \setminus \mathcal{B}^+$	$[\varepsilon]$	$[0, \bar{\varepsilon}_1(1 - \bar{\varepsilon}_2)]$	$\mathcal{A}^- \setminus \mathcal{B}^-$	$[\varepsilon]$	$[0, \bar{\omega}_2(1 - \underline{\omega}_1)]$
	$[\omega]$	$[\varepsilon_2]$		$[\omega]$	$[\omega_1]$
$\overline{\mathcal{A}^+}$	$[\omega]$	$[\varepsilon_1]$	$\overline{\mathcal{A}^-}$	$[\varepsilon]$	$[\omega_1]$

8 Data types that model random approximate sets

A *data type* is a set and the elements of the set are called the *values* of the data type. We impose a *structure* on sets (data types) by defining morphisms between them, such as operations like *intersection* or relations like *subset*. Morphisms are also types. Any data type needs one or more *value constructors*, functions that map to values of the type.

The random approximate set is an abstract data type that models a *set* with an additional set of *probabilistic* axioms described in section 3. Suppose T is a data type that overloads the *member-of* predicate $\in: \mathcal{U} \times T \mapsto \{0, 1\}$ and has a *value constructor* $\varepsilon\tau$ that is a *conditional probability distribution* over values of T given elements of type $2^{\mathcal{U}}$. Data type T models the abstract data type of the random approximate set over elements in \mathcal{U} with a false positive rate ε and true positive rate τ if axioms 2 and 3 are satisfied, i.e.,

$$\mathbb{P}[x \in \varepsilon\tau(\mathcal{S}) \mid x \notin \mathcal{S}] = \varepsilon \quad (8.1)$$

and

$$\mathbb{P}[x \in \varepsilon\tau(\mathcal{S}) \mid x \in \mathcal{S}] = \tau. \quad (8.2)$$

An instance of T also models a classic set by its membership predicate, i.e., two sets are *equal* if and only if they have the same members. We denote that an instance of T models a set \mathcal{A} by $T(\mathcal{A})$.

Normally, two different data types that model an abstract data type are *exchangable* over a set of *regular functions* without changing the result. However, random approximate sets are *probabilistic* so this strict definition of exchangability does not capture the intended meaning. The random approximate set model is a *frequentistic probability* model where an event's probability is defined as the *limit* of its relative frequency in a large number of trials. Thus, we relax the definition of exchangability and conclude that two data types that model random approximate sets (or any other probabilistic abstract data type) should produce the same *limit* of the relative frequency of results in a large number of *independent runs*.

An important distinction must be made with respect to *independent runs*. The most straightforward meaning is, given any set $\mathcal{A} \in 2^{\mathcal{U}}$, at the limit, repeated applications of $\varepsilon\tau(\mathcal{A})$ generates a sample that converges in distribution to $\mathcal{A}_\varepsilon^\tau$. However, we also wish to allow for *deterministic* value constructors.⁴

8.1 Deterministic value constructors

Value constructors compatible with the random approximate set model may come in many forms. For example, in section 9 we demonstrate an approximation of Boolean search where Boolean queries are deterministically mapped to approximate result sets compatible with the random approximate set model.

Suppose $\varepsilon\varepsilon: 2^{\mathcal{U}} \mapsto T$ (i.e., a deterministic total function) maps sets in $2^{\mathcal{U}}$ to objects of type T that model random approximate sets over the input. Since T models the abstract data type of the set, there is a unique bijection between T and $2^{\mathcal{U}}$, i.e., every value in T models a specific subset of \mathcal{U} . Thus, we may view $\varepsilon\varepsilon$ as a function $\varepsilon\tau: 2^{\mathcal{U}} \mapsto 2^{\mathcal{U}}$ with an *image*

$$\text{image}(\varepsilon\tau) = \{ \varepsilon\tau(\mathcal{A}) \mid \mathcal{A} \in 2^{\mathcal{U}} \} \subseteq 2^{\mathcal{U}}. \quad (8.3)$$

Since the value constructor $\varepsilon\tau$ may map multiple input sets to the same output set and some sets in the codomain may not be mapped to by any set in the domain, $\varepsilon\tau$ is (possibly) a non-surjective, non-injective function.

Definition 8.1. A σ -algebra is closed under countable unions, intersections, and complements.

The image of $\varepsilon\tau$ 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 along with the empty set \emptyset and the universal set \mathcal{U} is by definition a σ -algebra and is denoted by $\sigma(\varepsilon\varepsilon)$.

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

$$\sigma(\varepsilon\varepsilon) \mapsto \sigma(\varepsilon'\varepsilon') \quad (8.4)$$

and

$$\sigma(\varepsilon\varepsilon) \times \sigma(\varepsilon\varepsilon) \mapsto \sigma(\varepsilon\varepsilon). \quad (8.5)$$

Suppose we have two Boolean algebras, $(\sigma(\varepsilon), \cup, \cap, \bar{}, \emptyset, \mathcal{U})$ and $(\sigma(\mathfrak{g}), \cup, \cap, \bar{}, \emptyset, \mathcal{U})$, where $\varepsilon\tau$ and \mathfrak{g} are value constructors for approximate sets over $2^{\mathcal{U}}$. Set-theoretic operations over both Boolean algebras is the Boolean algebra $(\Sigma(\varepsilon\tau, \mathfrak{g}), \cup, \cap, \bar{}, \emptyset, \mathcal{U})$ where $\Sigma(\varepsilon\tau, \mathfrak{g}) = \sigma(\sigma(\varepsilon) \cup \sigma(\mathfrak{g}))$. Note that $\sigma(\varepsilon\tau), \sigma(\mathfrak{g}) \subseteq \Sigma(\varepsilon\tau, \mathfrak{g})$, so $\Sigma(f_1, \dots, f_n)$ converges to $2^{\mathcal{U}}$ as $n \rightarrow \infty$, where f_1, \dots, f_n are different mappings.

⁴Deterministic algorithms compatible with the random approximate set model are common but frequently have an auxiliary seed which indexes a particular approximation in a family.

Remark. It is often trivial to implement a family of deterministic value constructors $2^{\mathcal{U}} \mapsto T = \{f_1, \dots, f_n\}$ with distinct σ -algebras where T models random approximate sets over $2^{\mathcal{U}}$. Additionally, assuming each time an approximate set is constructed, a “random” value constructor from $2^{\mathcal{U}} \mapsto T$ is invoked, then repeated invocations on some set $\mathcal{A} \in 2^{\mathcal{U}}$ generates a frequency distribution of sets that converges to \mathcal{A}^\pm as $n \rightarrow \infty$, e.g., “randomly” seeding a Bloom filter’s hash function. \triangle

How do we reconcile a deterministic value constructor $\varepsilon\tau: 2^{\mathcal{U}} \mapsto T$ with the *probabilistic model*? In this context, the notion of *probability* quantifies our *ignorance*:

1. Given a set \mathcal{S} , we do not have complete *a priori* knowledge about the set the value constructor maps to. The approximate set model only provides *a priori* knowledge about the probability distribution \mathcal{S}^\pm . We acquire *a posteriori* knowledge⁵ by observing $\varepsilon\tau(\mathcal{S})$.
2. Given $T(\mathcal{S})$, 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 evaluating $\varepsilon\tau(\mathcal{A})$ for each $\mathcal{A} \in 2^{\mathcal{U}}$ and remembering the sets that map to $T(\mathcal{S})$.⁶ However, since $\varepsilon\tau$ is (possibly) non-injective, one or more sets may map to $T(\mathcal{S})$ and thus this process may not completely eliminate uncertainty. Additionally, the domain $2^{\mathcal{U}}$ has a cardinality $2^{|\mathcal{U}|}$ and thus exhaustive searches are impractical to compute even for relatively small domains.⁷

Suppose \mathcal{U} is finite. The set of deterministic value constructors $2^{\mathcal{U}} \mapsto 2^{\mathcal{U}}$ has a cardinality $(2^u)^{(2^u)}$, and in a sense they are all compatible with the random approximate set 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 $|U| = 2$, 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

$$\hat{\varepsilon}(\mathcal{X}) = \frac{1}{n} \sum_{x \in \bar{\mathcal{X}}} \mathbb{1}_{\varepsilon\tau(\mathcal{X})}(x), \tag{8.6}$$

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

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

$$\bar{\varepsilon} = \frac{1}{|\mathcal{U}_p|} \sum_{\mathcal{X} \in \mathcal{U}_p} \hat{\varepsilon}(\mathcal{X}), \tag{8.7}$$

is an unbiased estimator of ε and the population variance

$$s_\varepsilon^2 = \frac{1}{|\mathcal{U}_p|} \sum_{\mathcal{X} \in \mathcal{U}_p} \hat{\varepsilon}(\mathcal{X}), \tag{8.8}$$

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

Proof. We imagine that the function $\varepsilon\tau$ 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

⁵A posteriori knowledge is dependent on experience.

⁶If the approximate set is the result of the union, intersection, and complement of two or more approximate sets, then we must consider the closure.

⁷In the case of *countably infinite* domains, it is not even theoretically possible.

realization of $\mathcal{E}_n = \hat{\varepsilon}$. Thus,

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

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

□

8.2 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

$$\mathcal{O}(u) \text{ bits} \quad (\text{8.9})$$

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

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

Postulate 8.1. *The information-theoretic lower-bound of a data structure that implements the countably infinite random approximate set abstract data type has an expected bit length given by*

$$- \tau \log_2 \varepsilon \text{ bits/element}, \quad (\text{8.10})$$

where $\varepsilon > 0$ is the false positive rate and τ is the true positive.

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 ,

$$\text{RE}(X, Y) = \frac{\ell(Y)}{\ell(X)}, \quad (\text{8.11})$$

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 8.2. *The absolute space efficiency of a data structure X , denoted by $\mathbf{E}(X)$, 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 ,*

$$\mathbf{E}(X) = \frac{\theta}{\ell(X)}, \quad (\text{8.12})$$

where $\ell(X)$ denotes the bit length of X and θ denotes the bit length of the information-theoretic lower-bound. The closer $\mathbf{E}(X)$ 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 a random approximate set of an objective set with p elements with a false positive rate ε and true positive rate τ is given by

$$\mathbf{E}(X) = \frac{-p\tau \log_2 \varepsilon}{\ell(X)}. \quad (\text{8.13})$$

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

$$- \frac{1}{\ln 2} \log_2 \varepsilon \text{ bits/element}, \quad (\text{8.14})$$

thus the absolute efficiency of the Bloom filter is $\ln 2 \approx 0.69$. A practical implementation of the *positive random approximate set* is given by the *Perfect Hash Filter*[], which compares favorably to the Bloom filter in many circumstances.⁹

In [], we claimed that the method of moments estimator for p of an objective set given a particular realization of a random approximation set is undefined for countably infinite universes. Suppose we have a data structure X that *models* random approximate sets with an *expected* space complexity proportional to p , i.e., $p \cdot b(\tau, \varepsilon)$ bits, where b is the expected bits per *positive* element given a false positive rate ε and true positive rate τ . Then, given an object x of type X , an estimator of p is

$$\hat{p} = \frac{\ell(x)}{b(\tau, \varepsilon)}, \quad (\text{8.15})$$

where ℓ is the bit length function. An expected *upper-bound* on the cardinality is obtained by plugging in the information-theoretic lower-bound $b(\tau, \varepsilon) = -\tau \log_2 \varepsilon$ bits per element.

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, the number of positives p being the most obvious parameter.

⁹The *Singular Hash Set*[] is an example of a data structure that obtains optimality using *brute-force* search, so it is not practical for even relatively small objective sets. However, its primary purpose is analytic tractability.

8.2.1 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 $\mathcal{S}_1 = \{x_{j_1}, \dots, x_{j_m}\}$ and suppose we wish to insert the elements x_{k_1}, \dots, x_{k_p} into \mathcal{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 $\mathcal{S}_2 = \{x_{k_1}, \dots, x_{k_p}\}$, and then apply the union function,

$$X_1 \cup X_2. \quad (8.16)$$

If we replace X_1 and X_2 by objects that implement *positive approximate sets* of \mathcal{S}_1 and \mathcal{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 8.1. *Given two countably infinite positive approximate sets \mathcal{S}_1^+ and \mathcal{S}_2^+ respectively with false positive rates $\hat{\epsilon}_1$ and $\hat{\epsilon}_2$, the approximate set $\mathcal{S}_1^+ \cup \mathcal{S}_2^+$, 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*

$$\mathbf{E}(\hat{\epsilon}_1, \hat{\epsilon}_2 | \alpha_1, \alpha_2) = \frac{\log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 - \hat{\epsilon}_1 \hat{\epsilon}_2)}{\alpha_1 \log_2 \hat{\epsilon}_1 + \alpha_2 \log_2 \hat{\epsilon}_2}, \quad (8.17)$$

where

$$\begin{aligned} 0 < \alpha_1 &= \frac{|\mathcal{S}_1|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \leq 1, \\ 0 < \alpha_2 &= \frac{|\mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \leq 1, \\ 1 &\leq \alpha_1 + \alpha_2. \end{aligned} \quad (8.18)$$

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. ¹⁰

Proof. The proof comes in two parts. First, we prove eq. (8.17), and then we prove the bounds on α_1 and α_2 given by eq. (8.18).

Let X and Y denote optimally space-efficient data structures that respectively implement positive approximate sets \mathcal{S}_1^+ and \mathcal{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 $\mathcal{S}_1 \cup \mathcal{S}_2$ with the above false positive rate is given by

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

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, $\ell(X) = -|\mathcal{S}_1| \log_2 \hat{\epsilon}_1$ and $\ell(Y) = -|\mathcal{S}_2| \log_2 \hat{\epsilon}_2$. Making these substitutions yields an absolute efficiency ???

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

Letting

$$\alpha_1 = \frac{|\mathcal{S}_1|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \text{ and } \alpha_2 = \frac{|\mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|}, \quad (d)$$

we may rewrite eq. (c) as

$$\frac{\log_2(\hat{\epsilon}_1 + \hat{\epsilon}_2 - \hat{\epsilon}_1 \hat{\epsilon}_2)}{\alpha_1 \log_2 \hat{\epsilon}_1 + \alpha_2 \log_2 \hat{\epsilon}_2}. \quad (8.17 \text{ revisited})$$

In the second part of the proof, we prove the bounds on α_1 and α_2 as given by eq. (8.18). Both α_1 and α_2 must be non-negative since each is the ratio of two positive numbers (cardinalities). If $|\mathcal{S}_1| \ll |\mathcal{S}_2|$, then $\alpha_1 \approx 0$. If $\mathcal{S}_1 \supset \mathcal{S}_2$, then $\alpha_1 = 1$. A similar argument holds for α_2 . Finally,

$$\alpha_1 + \alpha_2 = \frac{|\mathcal{S}_1| + |\mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \quad (e)$$

has a minimum value by assuming that \mathcal{S}_1 and \mathcal{S}_2 are disjoint sets (i.e., their intersection is the empty set), in which case

$$\alpha_1 + \alpha_2 = \frac{|\mathcal{S}_1| + |\mathcal{S}_2|}{|\mathcal{S}_1| + |\mathcal{S}_2|} = 1. \quad (f)$$

□

¹⁰As $(\hat{\epsilon}_1, \hat{\epsilon}_2) \rightarrow (0, 0)$, the absolute efficiency depends on the path taken. In most cases, it goes to 0.

Figure 4: Expected lower-bound on efficiency of the union of two approximate sets, neither of which can be enumerated.

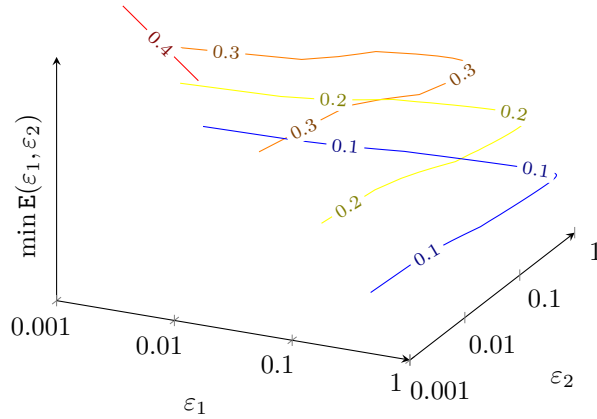
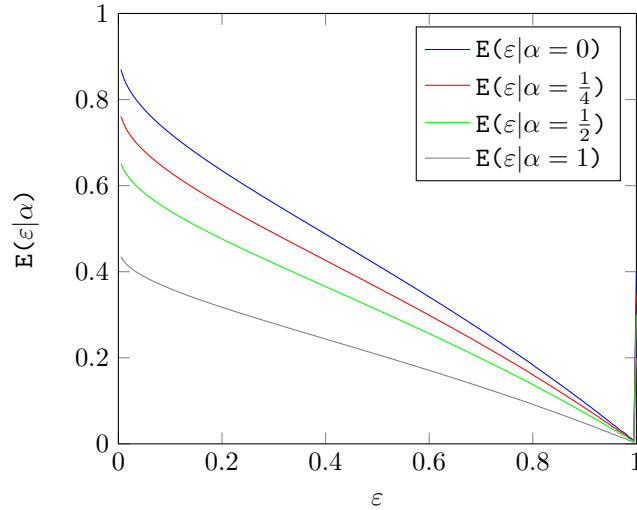


Figure 5: Expected lower-bound on efficiency of the union of two approximate sets with the same false positive rate ε , neither of which can be enumerated.



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

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

Corollary 8.1.1. *Given two positive, non-enumerable approximate sets with false positive rates $\hat{\varepsilon}_1$ and $\hat{\varepsilon}_2$, their union is an approximate set that has an efficiency that is expected to be greater than the lower bound given by*

$$\min \mathbf{E}(\hat{\varepsilon}_1, \hat{\varepsilon}_2) = \frac{\log_2(\hat{\varepsilon}_1 + \hat{\varepsilon}_2 - \hat{\varepsilon}_1 \hat{\varepsilon}_2)}{\log_2 \hat{\varepsilon}_1 \hat{\varepsilon}_2}. \quad (8.19)$$

Corollary 8.1.2. *If $\varepsilon_1 = \varepsilon_2 = \varepsilon$, then the absolute efficiency is given by*

$$\mathbf{E}(\hat{\varepsilon}|\alpha) = \left(1 + \frac{\log_2(2 - \hat{\varepsilon})}{\log_2 \hat{\varepsilon}}\right) \left(1 - \frac{\alpha}{2}\right), \quad (8.20)$$

where

$$0 \leq \alpha = \frac{|\mathcal{S}_1 \cap \mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \leq 1, \quad (8.21)$$

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

See fig. 5 for a graphic illustration.

8.3 code tmp

8.4 Algebra of sets

The perfect hash filter is a *first-order* rate-distorted set, which is a type of *random approximate set* whose error rates are due to rate-distortion.

Applying binary operators like *union* or *intersection* map pairs of first-order random approximate sets to *second-order* random approximate sets. If we continue in this trend, we generate higher-order random approximate sets, e.g., the union of a first-order set and a second-order set is a third-order set.

Note, however, that *complements* of an n -th order set is an n -order set, i.e., the order of approximation is closed under complements.

When take the union of a pair of first-order random approximate sets, the result is a second-order set.

```

c++ template < FirstOrderApproximateSet A, FirstOrderApproximateSet B > struct second_order_approximate_set_union_expr_using_g...
auto contains(value_type x) return a.contains(x) || b.contains(x); A const; B const;
c++ auto fpr(second_order_approximate_set_union_expr A) // composed sets should be points. we decide to // span the points yielding a single
c++ template < typename X, FirstOrderApproximateSet A, FirstOrderApproximateSet B > SecondOrderApproximate-
SetUnionExpr<X> operator+(A a, B b) return SecondOrderApproximateSetUnionExpr<X>(a,b);
c++ template < typename X, FirstOrderApproximateSet A, FirstOrderApproximateSet B > SecondOrderApproximate-
Set<X> operator*(ApproximateSet<X> a, ApproximateSet<X> b) return ( a + b);

```

Using type-erasure, we wrap arbitrarily complex expressions into `c++ ApproximateSet<X>`. `c++ template < typename X, FirstOrderApproximateSet A, FirstOrderApproximateSet B > ApproximateSet<X> operator*(ApproximateSet<X> a, ApproximateSet<X> b) return (a + b);`

9 Application: approximating Boolean search

An information retrieval process begins when a user submits a *query* to an information system, where a query represents an *information need*. In response, the information system returns a set of relevant documents that satisfy the query.

Boolean search is an information retrieval model given by the following definition.

Definition 9.1. *A document in the collection is either relevant or non-relevant to a Boolean query.*

We do not specify the structure of documents since we are only interested in being able to specify documents in a collection by some *label*, e.g., a file name. We specify the universal set of document labels by \mathcal{D} and therefore a *particular* collection of interest is a subset of \mathcal{D} .

We consider queries over the Boolean algebra $Q = (2^{\mathcal{K}}, \wedge, \vee, \neg, \epsilon, \mathcal{K})$, where \mathcal{K} denotes a set of *search keys*, e.g., units of information like English words.¹¹ Without loss of generality, we transform Boolean queries over Q to the BNF

$$\begin{aligned}
\langle \text{query} \rangle &:= \text{“}\langle \text{key} \rangle\text{”} \mid \neg (\langle \text{query} \rangle) \mid \\
&\quad \vee (\langle \text{query} \rangle, \langle \text{query} \rangle) \mid \\
&\quad \wedge (\langle \text{query} \rangle, \langle \text{query} \rangle) \\
\langle \text{key} \rangle &:= \text{a key in } \mathcal{K}.
\end{aligned}$$

Search indexes may be used to quickly compute whether a given document is relevant to a given query. Since we are using a Boolean search query model Q , search indexes may be efficiently represented by *sets* over \mathcal{K} in the Boolean algebra $S = (2^{\mathcal{K}}, \cap, \cup, \bar{}, \emptyset, \mathcal{K})$. In particular, let $d : \mathcal{D} \mapsto 2^{\mathcal{K}}$ be a function that maps documents to search indexes with a definition given by

$$(d) := \{ k \in \mathcal{K} \mid k \text{ is relevant to } d \}. \quad (9.1)$$

The set of relevant documents to a query is denoted the query’s *result set*. The result sets form the Boolean algebra $R = (2^{\mathcal{D}}, \cap, \cup, \bar{}, \emptyset, \mathcal{D})$.

A bijection from Q to S is given by $\wedge \mapsto \cap, \vee \mapsto \cup, \neg \mapsto \bar{}, \epsilon \mapsto \emptyset$, and $\mathcal{K} \mapsto \mathcal{K}$. Let $\text{find} : Q \times 2^{\mathcal{D}} \mapsto 2^{\mathcal{D}}$ be the function that maps queries in Q to result sets in R by using the collection of corresponding search indexes in S ,

$$\text{find}(q, ds) := \begin{cases} \bar{}(\text{find}(t, ds)) & \text{if } h = \neg \\ \cup(\text{find}(\text{left}(t), ds), \text{find}(\text{right}(t), ds)) & \text{if } h = \vee \\ \cap(\text{find}(\text{left}(t), ds), \text{find}(\text{right}(t), ds)) & \text{if } h = \wedge \\ \{ d \in ds \mid h \in (d) \} & \text{otherwise,} \end{cases} \quad (9.2)$$

where $h = \text{head}(q)$, $t = \text{tail}(q)$, $\text{head} : Q \mapsto \{\neg, \vee, \wedge\} \cup \mathcal{K}$ maps any given query q to the next *Boolean* operation or *key* in q , $\text{tail} : Q \mapsto Q$ maps any given query q to nested sub-queries in q , e.g., $\text{tail}(\vee(q_1, q_2)) = ?$, left maps $f(x, y)$ to x and right maps $f(x, y)$ to y .

¹¹This is isomorphic to the Boolean algebra $(\{0, 1\}^k, \wedge, \vee, \neg, 0^k, 1^k)$ where $k = |\mathcal{K}|$.

9.1 Random approximate Boolean search

We consider an *approximation* of the set-theoretic *Boolean search* model where the Boolean search indexes are replaced by *random approximate sets*, i.e., $\cdot : \mathcal{D} \mapsto 2^{\mathcal{K}}$ is replaced with $\sigma : \mathcal{D} \mapsto 2^{\mathcal{K}^\pm}$. We denote the transformed find function by \mathbf{find}^σ as opposed to the objective function \mathbf{find} .

This replacement *induces* random approximate *result sets* as given by the following theorem.

Theorem 9.1. \mathbf{find}^σ is an approximation of \mathbf{find} where $\mathbf{find}^\sigma(q, ds)$ is a random approximate set of $\mathbf{find}(q, ds)$ for all q in Q and all $ds \in 2^{\mathcal{D}}$.

Proof. An approximate search index $\mathcal{S}_\omega^\varepsilon$ is relevant to a key x if the key x tests positive in it. A false positive occurs if the key x is not in \mathcal{S} but is in \mathcal{S}^\pm , which occurs with some probability $\varepsilon \geq 0$. A false negative occurs if a key x is in \mathcal{S} but is not in \mathcal{S}^\pm , which occurs with some probability $\omega \geq 0$.

We have established that the result sets of a single atomic key are approximate result sets. We may now apply the set-theoretic results in ?? to implement the full set-theoretic model for approximate sets.

Continue proof here. □

\mathbf{find}^σ is a function and therefore produces the same output (result set) for the same input (query). However, it still obeys the axioms of the random approximate set model since as described in ??.

Example 6 Suppose the search indexes are positive approximate sets each with a false positive rate ε . A common type of Boolean query is the intersection (conjunction) of atomic keys. Consider a conjunctive query of k keys, x_1, \dots, x_k . The result set $\mathcal{R}^+ (\{x_1\} \cap \dots \cap \{x_k\}) = \mathbf{find}(\overline{x_1} \cup \dots \cup \overline{x_k})$ is a positive approximate set with an uncertain false positive rate $[\varepsilon_k] = [\varepsilon^k, \varepsilon]$.

Proof. Let the approximate result set for key x_j be denoted by $\mathcal{R}_{x_j}^+$. The result set is given by

$$\mathcal{R}_{\mathcal{X}}^+ = \bigcap_{j=1}^k \mathcal{R}_{x_j}^+. \quad (\text{a})$$

By ??, \mathcal{R}_j^+ has a false positive rate ε for $j \in [1, \dots, k]$. By ??, $\mathcal{R}_1^+ \cap \mathcal{R}_2^+$ has a false positive rate $[\varepsilon^2, \varepsilon]$. Similarly,

$$\left(\mathcal{R}_1^+ \cap \mathcal{R}_2^+ \right) \cap \mathcal{R}_3^+ = \mathcal{R}_1^+ \cap \mathcal{R}_2^+ \cap \mathcal{R}_3^+ \quad (\text{b})$$

has a false positive rate $[\varepsilon^3, \varepsilon]$. Continuing in this fashion, we see that $\mathcal{R}_{\mathcal{X}}^+ = \mathcal{R}_1^+ \cap \dots \cap \mathcal{R}_k^+$ has a false positive rate $[\varepsilon^k, \varepsilon]$. □

To quantify the performance measure of the information retrieval system, we may use the binary classification results in section 6.

A Proof of corollary 4.2.2

To say that the sequence $\mathcal{E}_1, \mathcal{E}_2, \dots$ converges almost surely to ε means that

$$\mathbb{P} \left[\lim_{n \rightarrow \infty} \mathcal{E}_n = \varepsilon \right] = 1. \quad (\text{A.1})$$

By corollary 4.2.2, given *countably infinite* negatives, a random approximate set with a false positive rate ε is *certain* to obtain ε .

Proof. Hoeffding's inequality[?] provides that FP_n is concentrated around its mean $n\varepsilon$ as given by

$$\mathbb{P}[(\varepsilon - \epsilon)n \leq \text{FP}_n \leq (\varepsilon + \epsilon)n] \geq 1 - 2 \exp(-2\epsilon^2 n), \quad (\text{a})$$

where $\epsilon > 0$. We are interested in the limiting probability

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{P}[(\varepsilon - \epsilon)n \leq \text{FP}_n \leq (\varepsilon + \epsilon)n] = \\ \lim_{n \rightarrow \infty} \left\{ 1 - 2 \exp(-2\epsilon^2 n) \right\} = 1. \end{aligned} \quad (\text{b})$$

As ϵ goes to 0, $\lim_{n \rightarrow \infty} \text{FP}_n$ converges almost surely to εn and therefore $\lim_{n \rightarrow \infty} \text{FP}_n/n$ converges almost surely to ε . □

B Proof of theorem 6.1

Given p positives and n negatives, by ?? an approximate set with a false positive rate ε and a false negative rate ω has an *expected* precision given *approximately* by

$$\text{ppv}(\omega, \varepsilon; n, p) \approx \frac{\bar{t}_p}{\bar{t}_p + \bar{f}_p} + \frac{\bar{t}_p \sigma_{f_p}^2 - \bar{f}_p \sigma_{t_p}^2}{(\bar{t}_p + \bar{f}_p)^3}, \quad (?? \text{ revisited})$$

where $\bar{t}_p = p\tau$ is the *expected* number of *true positives*, $\bar{f}_p = n\varepsilon$ is the *expected* number of *false positives*, $\sigma_{t_p}^2 = p\omega\tau$ is the variance of the number of *true positives*, and $\sigma_{f_p}^2 = n\varepsilon\omega$ is the variance of the number of *false positives*.

Proof. The positive predictive value is a random variable given by

$$\frac{\text{TP}_p}{\text{TP}_p + \text{FP}_n}. \quad (\text{a})$$

We are interested in the *expected* positive predictive value,

$$\text{ppv}(\varepsilon, \tau) = \text{E} \left[\frac{\text{TP}_p}{\text{TP}_p + \text{FP}_n} \right]. \quad (\text{b})$$

This expectation is of a non-linear function of random variables, which is problematic so we choose to approximate the expectation.

Let the *positive predictive value* function be denoted by

$$f(t_p, f_p) = \frac{t_p}{t_p + f_p}, \quad (\text{c})$$

where t_p is the number of true positives and f_p is the number of false positives. We approximate this function with a second-order Taylor series. The gradient of f is given by

$$\nabla f(t_p, f_p) = \frac{1}{(t_p + f_p)^2} \begin{bmatrix} f_p \\ -t_p \end{bmatrix} \quad (\text{d})$$

and the Hessian of f is given by

$$\mathcal{H}(t_p, f_p) = \frac{1}{(t_p + f_p)^3} \begin{bmatrix} -2f_p & t_p - f_p \\ t_p - f_p & 2t_p \end{bmatrix}. \quad (\text{e})$$

A linear approximation g of f that is reasonably accurate near the expected value of TP_p , denoted by \bar{t}_p , and the expected value of FP_n , denoted by \bar{f}_p , is given by

$$g(t_p, f_p) = f(\bar{t}_p, \bar{f}_p) + \nabla f(\bar{t}_p, \bar{f}_p)^\top \begin{bmatrix} t_p - \bar{t}_p \\ f_p - \bar{f}_p \end{bmatrix} + \frac{1}{2} \begin{bmatrix} t_p - \bar{t}_p \\ f_p - \bar{f}_p \end{bmatrix}^\top \mathcal{H}(\bar{t}_p, \bar{f}_p) \begin{bmatrix} t_p - \bar{t}_p \\ f_p - \bar{f}_p \end{bmatrix}. \quad (\text{f})$$

As a function of random variables TP_p and FP_n , $g(\text{TP}_p, \text{FP}_n)$ is a random variable. Since $\text{E}[\text{TP}_p - \bar{t}_p] = 0$ and $\text{E}[\text{FP}_n - \bar{f}_p] = 0$, we immediately simplify the expectation of g to

$$\text{E}[g(\text{TP}_p, \text{FP}_n)] = \frac{\bar{t}_p}{\bar{t}_p + \bar{f}_p} + \frac{\text{E}[A]}{(\bar{t}_p + \bar{f}_p)^3} \quad (\text{g})$$

where

$$A = \frac{1}{2} \begin{bmatrix} \text{TP}_p - \bar{t}_p \\ \text{FP}_n - \bar{f}_p \end{bmatrix}^\top \begin{bmatrix} -2\bar{f}_p & \bar{t}_p - \bar{f}_p \\ \bar{t}_p - \bar{f}_p & 2\bar{t}_p \end{bmatrix} \begin{bmatrix} \text{TP}_p - \bar{t}_p \\ \text{FP}_n - \bar{f}_p \end{bmatrix}. \quad (\text{h})$$

Multiplying the right column matrix by the Hessian matrix in A results in

$$A = \frac{1}{2} \begin{bmatrix} \text{TP}_p - \bar{t}_p \\ \text{FP}_n - \bar{f}_p \end{bmatrix}^\top \begin{bmatrix} -2\bar{f}(\text{TP}_p - \bar{t}_p) + (\bar{t}_p - \bar{f}_p)(\text{FP}_n - \bar{f}_p) \\ (\bar{t}_p - \bar{f}_p)(\text{TP}_p - \bar{t}_p) + 2\bar{t}_p(\text{FP}_n - \bar{f}_p) \end{bmatrix} \quad (\text{i})$$

Multiplying the left column matrix by the right column matrix in A results in

$$A = -\bar{f}_p (\text{TP}_p - \bar{t}_p)^2 + (\bar{t}_p - \bar{f}_p) (\text{TP}_p - \bar{t}_p) (\text{FP}_n - \bar{f}_p) + \bar{t}_p (\text{FP}_n - \bar{f}_p)^2. \quad (\text{j})$$

As a linear operator, the expectation of A is equivalent to

$$\mathbb{E}[A] = -\bar{f}_p \mathbb{E}[\text{TP}_p - \bar{t}_p]^2 + (\bar{t} - \bar{f}_p) \mathbb{E}\left[\left(\text{FP}_n - \bar{f}_p\right) (\text{TP}_p - \bar{t}_p)\right] + \bar{t}_p \mathbb{E}\left[\text{FP}_n - \bar{f}_p\right]^2. \quad (\text{k})$$

By definition, $\mathbb{E}[\text{TP}_p - \bar{t}_p]^2$ is the variance of TP_p , $\mathbb{E}\left[\left(\text{FP}_n - \bar{f}_p\right)^2\right]$ is the variance of FP_n , and $\mathbb{E}\left[\left(\text{FP}_n - \bar{f}_p\right) (\text{TP}_p - \bar{t}_p)\right]$ is the covariance of TP_p and FP_n , which is 0 since they are independent. Making these substitutions results in

$$\mathbb{E}[A] = \bar{t}_p \text{V}[\text{FP}_n] - \bar{f}_p \text{V}[\text{TP}_p]. \quad (\text{l})$$

Substituting this result into eq. (g) yields

$$\mathbb{E}[\text{g}(\text{TP}_p, \text{FP}_n)] = \frac{\bar{t}_p}{\bar{t}_p + \bar{f}_p} + \frac{-\bar{f}_p \text{V}[\text{TP}_p] + \bar{t}_p \text{V}[\text{FP}_n]}{(\bar{t}_p + \bar{f}_p)^3} \quad (\text{m})$$

By theorem 4.1, FP_n is binomially distributed with a mean $n\varepsilon$ and a variance $n\varepsilon\eta$ and by corollary 4.4.2, TP_p is binomially distributed with a mean $p\tau$ and a variance $p\omega\tau$. \square

C Proof of theorem 5.3

Given $\mathcal{A}^\pm(\varepsilon_1, \omega)$ and $\mathcal{B}^\pm(\varepsilon_2, \omega_2)$, their union is an *approximate set* with a false negative rate given by

$$\begin{aligned} \omega &= \alpha_1 \omega_1 (1 - \varepsilon_2) + \alpha_2 \omega_2 (1 - \varepsilon_1) \\ &\quad + (1 - \alpha_1 - \alpha_2) \omega_1 \omega_2, \end{aligned} \quad (\text{5.6 revisited})$$

where

$$\begin{aligned} 0 \leq \alpha_1 &= \frac{|\mathcal{S}_1 \setminus \mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|}, \\ 0 \leq \alpha_2 &= \frac{|\mathcal{S}_2 \setminus \mathcal{S}_1|}{|\mathcal{S}_1 \cup \mathcal{S}_2|}, \\ \alpha_1 + \alpha_2 &\leq 1. \end{aligned} \quad (\text{5.7 revisited})$$

Proof. Suppose we have two sets \mathcal{S}_1 and \mathcal{S}_2 . The false negative rate ω is a probability conditioned on a positive in the union of sets \mathcal{S}_1 and \mathcal{S}_2 being a negative in the union of approximate sets \mathcal{S}^1 and \mathcal{S}^2 .

The set of *possible* false negatives is the set of positives, $\mathcal{S}_1 \cup \mathcal{S}_2$, which is equivalent to the union of the disjoint sets $\mathcal{S}_1 \cap \mathcal{S}_2$, $\mathcal{S}_1 \setminus \mathcal{S}_2$ and $\mathcal{S}_2 \setminus \mathcal{S}_1$.

The false negative rate is equivalent to ratio of the *expected* number of false negatives to the maximum possible false negatives $|\mathcal{S}_1 \cup \mathcal{S}_2|$. Since they are disjoint, we may consider each independently to calculate the expected total number of false negatives.

Let A_1 denote the event $X \in \mathcal{S}^1$, A_2 denote $X \in \mathcal{S}^2$, B_1 denote $X \in \mathcal{S}_1$, and B_2 denote $X \in \mathcal{S}_2$. Suppose we randomly select an element from $\mathcal{S}_1 \cap \mathcal{S}_2$. The probability that X is a negative in $\mathcal{S}^1 \cup \mathcal{S}^2$ given that it is positive in $\mathcal{S}_1 \cap \mathcal{S}_2$ is given by

$$\omega_{1\cap 2} = \mathbb{P}[(A_1 \cup A_2)' \mid B_1 \cap B_2]. \quad (\text{a})$$

By De Morgan's law, $(A_1 \cup A_2)' \equiv A_1' \cap A_2'$. Making this substitution results in

$$\omega_{1\cap 2} = \mathbb{P}[A_1' \cap A_2' \mid B_1 \cap B_2]. \quad (\text{b})$$

Since A_1 and A_2 are independent events, by the rules of probability

$$\omega_{1\cap 2} = \mathbb{P}[A_1' \mid B_1 \cap B_2] \mathbb{P}[A_2' \mid B_1 \cap B_2]. \quad (\text{c})$$

Since A_1 is independent of B_2 and A_2 is independent of B_1 , by the rules of probability

$$\omega_{1\cap 2} = \mathbb{P}[A_1' \mid B_1] \mathbb{P}[A_2' \mid B_2]. \quad (\text{d})$$

By definition, $\mathbb{P}[A_j' \mid B_j]$ is the false negative rate ω_j . Making this substitution yields

$$\omega_{1\cap 2} = \omega_1 \omega_2. \quad (\text{e})$$

There are $|\mathcal{S}_1 \cap \mathcal{S}_2|$ elements in $\mathcal{S}_1 \cap \mathcal{S}_2$, where each is an independent Bernoulli trial. Thus, there are expected to be

$$|\mathcal{S}_1 \cap \mathcal{S}_2| \omega_{1 \cap 2} = |\mathcal{S}_1 \cap \mathcal{S}_2| \omega_1 \omega_2 \quad (\text{f})$$

false negatives in $\mathcal{S}_1 \cap \mathcal{S}_2$.

Suppose we randomly select an element from $\mathcal{S}_1 \setminus \mathcal{S}_2$. The probability that X is a negative in $\mathcal{S}^1 \cup \mathcal{S}^2$ given that it is a positive in $\mathcal{S}_1 \setminus \mathcal{S}_2$ is given by

$$\omega_{1 \cap \bar{2}} = \text{P}[A'_1 \cap A'_2 \mid B_1 \cap B'_2]. \quad (\text{g})$$

Since A_1 and A_2 are independent events, this may be rewritten as

$$\omega_{1 \cap \bar{2}} = \text{P}[A'_1 \mid B_1 \cap B'_2] \text{P}[A'_2 \mid B_1 \cap B'_2]. \quad (\text{h})$$

Since A_1 is independent of B_2 and A_2 is independent of B_1 , this may be rewritten as

$$\omega_{1 \cap \bar{2}} = \text{P}[A'_1 \mid B_1] \text{P}[A'_2 \mid B'_2]. \quad (\text{i})$$

By definition, $\text{P}[A'_1 \mid B_1]$ is the false negative rate ω_1 and $\text{P}[A'_2 \mid B'_2]$ is the false positive rate ε_2 . Thus,

$$\omega_{1 \cap \bar{2}} = \omega_1(1 - \varepsilon_2). \quad (\text{j})$$

There are $|\mathcal{S}_1 \setminus \mathcal{S}_2|$ elements in $\mathcal{S}_1 \setminus \mathcal{S}_2$, where each is an independent Bernoulli trial. Thus, there are expected to be

$$|\mathcal{S}_1 \setminus \mathcal{S}_2| \omega_{1 \cap \bar{2}} = |\mathcal{S}_1 \setminus \mathcal{S}_2| \omega_1(1 - \varepsilon_2) \quad (\text{k})$$

false negatives in $\mathcal{S}_1 \setminus \mathcal{S}_2$. A similar argument follows for $\mathcal{S}_2 \setminus \mathcal{S}_1$ where there are expected to be

$$|\mathcal{S}_2 \setminus \mathcal{S}_1| \omega_2(1 - \varepsilon_2) \quad (\text{l})$$

false negatives.

The false negative rate is given by the ratio of the total expected number of false negatives given by eqs. (f), (k) and (l) to the total number of possible false negatives $|\mathcal{S}_1 \cup \mathcal{S}_2|$, which is given by

$$\omega = \frac{|\mathcal{S}_1 \setminus \mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \omega_1(1 - \varepsilon_2) + \frac{|\mathcal{S}_2 \setminus \mathcal{S}_1|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \omega_2(1 - \varepsilon_1) + \frac{|\mathcal{S}_1 \cap \mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \omega_1 \omega_2. \quad (\text{m})$$

If we let

$$\alpha_1 = \frac{|\mathcal{S}_1 \setminus \mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|} \text{ and } \alpha_2 = \frac{|\mathcal{S}_2 \setminus \mathcal{S}_1|}{|\mathcal{S}_1 \cup \mathcal{S}_2|}, \quad (\text{n})$$

then

$$1 - \alpha_1 - \alpha_2 = \frac{|\mathcal{S}_1 \cap \mathcal{S}_2|}{|\mathcal{S}_1 \cup \mathcal{S}_2|}. \quad (\text{o})$$

Making these substitutions into eq. (m) yields the result

$$\omega = \alpha_1 \omega_1(1 - \varepsilon_2) + \alpha_2 \omega_2(1 - \varepsilon_1) + (1 - \alpha_1 - \alpha_2) \omega_1 \omega_2. \quad (\text{p})$$

□

D Sampling distribution of arbitrary functions

TODO: add generative model as an algorithm for approximate sets? Add C++ implementation of the model? Do some simulations to see how rapidly it converges to the normal? TODO: feed in something like ppv function and see how well it matches the solution given in that one section. etc.

Suppose we have an objective function $f: 2^{\mathcal{X}^1} \times \dots \times 2^{\mathcal{X}^q} \mapsto \mathcal{Y}$, and we are interested in evaluating the loss when we replace one or more of the objective input sets with particular corresponding random approximate sets. The result of this substitution, as previously described, is a probability distribution over \mathcal{Y} .

The probability distribution of random approximate sets are precisely given; therefore, we may estimate the distribution of any function of random approximate sets by generating the random approximate sets and applying the function of interest.

Consider the m -by- q matrix where the (i, j) -th element is the random approximate set $\mathcal{A}^{i,j}(\tau_j, \varepsilon_j)$ such that they are all independently distributed and $\mathcal{A}^{i,j}$ for $i = 1, \dots, m$ are also identically distributed. If we apply g to each row of the matrix,

$$Y_i = g(\mathcal{A}^{i,1}, \dots, \mathcal{A}^{i,q}) \quad (\text{D.1})$$

for $i = 1, \dots, m$, we generate m i.i.d. random elements Y_1, \dots, Y_m .

If \mathcal{Y} is a measure space (discrete or continuous), consider the random variable

$$\bar{Y}_m = \frac{1}{m} \sum_{i=1}^m Y_i. \tag{D.2}$$

If Y_1 has a well-defined mean and variance, then by the central limit theorem

$$\lim_{m \rightarrow \infty} \bar{Y}_m \tag{D.3}$$

converges in distribution to a normal with a mean $E[Y_1]$ and a variance $V[Y_1]/m$.

A general approach to estimating \bar{Y}_m is given by generating a large sample of matrices and applying the function g to each to generate a large sample from Y_1 .

We provide an implementation of the generative model and a tool set that permits one to analyze various properties of the distribution of the function of interest.