

# Set Similarity Search Beyond MinHash\*

Tobias Christiani  
IT University of Copenhagen  
Copenhagen, Denmark  
tobc@itu.dk

Rasmus Pagh  
IT University of Copenhagen  
Copenhagen, Denmark  
pagh@itu.dk

## ABSTRACT

We consider the problem of approximate set similarity search under Braun-Blanquet similarity  $B(x, y) = |x \cap y|/\max(|x|, |y|)$ . The  $(b_1, b_2)$ -approximate Braun-Blanquet similarity search problem is to preprocess a collection of sets  $P$  such that, given a query set  $q$ , if there exists  $x \in P$  with  $B(q, x) \geq b_1$ , then we can efficiently return  $x' \in P$  with  $B(q, x') > b_2$ .

We present a simple data structure that solves this problem with space usage  $O(n^{1+\rho} \log n + \sum_{x \in P} |x|)$  and query time  $O(|q|n^\rho \log n)$  where  $n = |P|$  and  $\rho = \log(1/b_1)/\log(1/b_2)$ . Making use of existing lower bounds for locality-sensitive hashing by O'Donnell et al. (TOCT 2014) we show that this value of  $\rho$  is tight across the parameter space, i.e., for every choice of constants  $0 < b_2 < b_1 < 1$ .

In the case where all sets have the same size our solution strictly improves upon the value of  $\rho$  that can be obtained through the use of state-of-the-art data-independent techniques in the Indyk-Motwani locality-sensitive hashing framework (STOC 1998) such as Broder's MinHash (CCS 1997) for Jaccard similarity and Andoni et al.'s cross-polytope LSH (NIPS 2015) for cosine similarity. Surprisingly, even though our solution is data-independent, for a large part of the parameter space we outperform the currently best data-dependent method by Andoni and Razenshteyn (STOC 2015).

## 1 INTRODUCTION

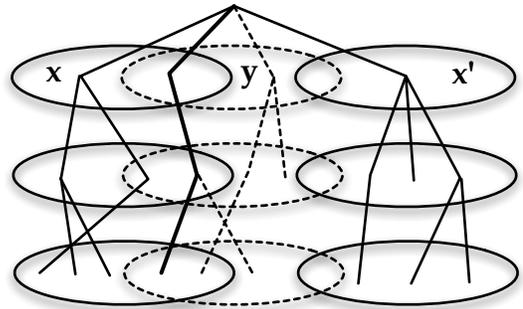
In this paper we consider the approximate set similarity problem or, equivalently, the problem of approximate Hamming near neighbor search in sparse vectors. Data that can be represented as sparse vectors is ubiquitous — a typical example is the representation of text documents as *term vectors*, where non-zero vector entries correspond to occurrences of words (or shingles). In order to perform identification of near-identical text documents in web-scale collections, Broder et al. [11, 12] designed and implemented *MinHash* (a.k.a. min-wise hashing), now understood as a locality-sensitive hash function [21]. This allowed approximate answers to similarity queries to be computed much faster than by other methods, and in particular made it possible to cluster the web pages of the AltaVista search engine (for the purpose of eliminating near-duplicate search results). Almost two decades after it was first described, MinHash remains one of the most widely used locality-sensitive hashing methods as witnessed by thousands of citations of [11, 12].

A *similarity measure* maps a pair of vectors to a similarity score in  $[0, 1]$ . It will often be convenient to interpret a vector  $x \in \{0, 1\}^d$  as the set  $\{i \mid x_i = 1\}$ . With this convention the *Jaccard similarity* of two vectors can be expressed as  $J(x, y) = |x \cap y|/|x \cup y|$ . In *approximate similarity search* we are interested the problem of

searching a data set  $P \subseteq \{0, 1\}^d$  for a vector of similarity at least  $j_1$  with a query vector  $q \in \{0, 1\}^d$ , but allow the search algorithm to return a vector of similarity  $j_2 < j_1$ . To simplify the exposition we will assume throughout the introduction that all vectors are *t*-sparse, i.e., have the same Hamming weight  $t$ .

Recent theoretical advances in data structures for approximate *near neighbor* search in Hamming space [5] make it possible to beat the asymptotic performance of MinHash-based Jaccard similarity search (using the LSH framework of [21]) in cases where the similarity threshold  $j_2$  is not too small. However, numerical computations suggest that MinHash is always better when  $j_2 < 1/45$ .

In this paper we address the problem: Can similarity search using MinHash be improved *in general*? We give an affirmative answer in the case where all sets have the same size  $t$  by introducing **CHOSEN PATH**: a simple data-independent search method that strictly improves MinHash, and is always better than the data-dependent method of [5] when  $j_2 < 1/9$ . Similar to data-independent locality-sensitive filtering (LSF) methods [9, 16, 24] our method works by mapping each data (or query) vector to a set of keys that must be stored (or looked up). The name **CHOSEN PATH** stems from the way the mapping is constructed: As paths in a layered random graph where the vertices at each layer is identified with the set  $\{1, \dots, d\}$  of dimensions, and where a vector  $x$  is only allowed to choose paths that stick to non-zero components  $x_i$ . This is illustrated in Figure 1.



**Figure 1:** CHOSEN PATH uses a branching process to associate each vector  $x \in \{0, 1\}^d$  with a set  $M_k(x) \subseteq \{1, \dots, d\}^k$  of paths of length  $k$  (in the picture  $k = 3$ ). The paths associated with  $x$  are limited to indices in the set  $\{i \mid x_i = 1\}$ , represented by an ellipsoid at each level in the illustration. In the example the set sizes are:  $|M_3(x)| = 4$  and  $|M_3(y)| = |M_3(x')| = 3$ . Parameters are chosen such that a query  $y$  that is similar to  $x \in P$  is likely to have a common path in  $x \cap y$  (shown as a bold line), whereas it shares few paths in expectation with vectors such as  $x'$  that are not similar.

\*The research leading to these results has received funding from the European Research Council under the European Union's 7th Framework Programme (FP7/2007-2013) / ERC grant agreement no. 614331.

## 1.1 Related Work

High-dimensional approximate similarity search methods can be characterized in terms of their  $\rho$ -value which is the exponent for which queries can be answered in time  $O(dn^\rho)$ , where  $n$  is the size of the set  $P$  and  $d$  denotes the dimensionality of the space. Here we focus on the “balanced” case where we aim for space  $O(n^{1+\rho} + dn)$ , but note that there now exist techniques for obtaining other trade-offs between query time and space overhead [4, 16].

**Locality-Sensitive Hashing Methods.** We begin by describing results for Hamming space, which is a special case of similarity search on the unit sphere (many of the results cited apply to the more general case). In Hamming space the focus has traditionally been on the  $\rho$ -value that can be obtained for solutions to the  $(r, cr)$ -approximate near neighbor problem: Preprocess a set of points  $P \subseteq \{0, 1\}^d$  such that, given a query point  $\mathbf{q}$ , if there exists  $\mathbf{x} \in P$  with  $\|\mathbf{x} - \mathbf{q}\|_1 \leq r$ , then return  $\mathbf{x}' \in P$  with  $\|\mathbf{x}' - \mathbf{q}\|_1 < cr$ . In the literature this problem is often presented as the  $c$ -approximate near neighbor problem where bounds for the  $\rho$ -value are stated in terms of  $c$  and, in the case of upper bounds, hold for every choice of  $r$ , while lower bounds may only hold for specific choices of  $r$ .

O’Donnell et al. [30] have shown that the value  $\rho = 1/c$  for  $c$ -approximate near neighbor search in Hamming space, obtained in the seminal work of Indyk and Motwani [23], is the best possible in terms of  $c$  for schemes based on Locality-Sensitive Hashing (LSH). However, the lower bound only applies when the distances of interest,  $r$  and  $cr$ , are relatively small compared to  $d$ , and better upper bounds are known for large distances. Notably, other LSH schemes for angular distance on the unit sphere such as cross-polytope LSH [2] give lower  $\rho$ -values for large distances. Extensions of the lower bound of [30] to cover more of the parameter space were recently given in [4, 16]. Until recently the best  $\rho$ -value known in terms of  $c$  was  $1/c$ , but in a sequence of papers Andoni et al. [3, 5] have shown how to use *data-dependent* LSH techniques to achieve  $\rho = 1/(2c-1) + o_n(1)$ , bypassing the lower bound framework of [30] which assumes the LSH to be independent of data.

**Set Similarity Search.** There exists a large number of different measures of set similarity with various applications for which it would be desirable to have efficient approximate similarity search algorithms [15]. Given a measure of similarity  $S$  assume that we have access to a family  $\mathcal{H}$  of locality-sensitive hash functions (defined in Section 2) such that for  $h \sim \mathcal{H}$  it holds for every pair of sets  $\mathbf{x}, \mathbf{y}$  that  $\Pr[h(\mathbf{x}) = h(\mathbf{y})] = S(\mathbf{x}, \mathbf{y})$ . Then we can use the LSH framework to construct a solution for the  $(s_1, s_2)$ -approximate similarity search problem under  $S$  with exponent  $\rho = \log(1/s_1)/\log(1/s_2)$ . With respect to the existence of such families Charikar [13] showed that if the similarity measure  $S$  admits an LSH with the above properties, then  $1 - S$  must be a metric. Recently, Chierichetti and Kumar [14] showed that, given a similarity  $S$  that admits an LSH with the above properties, the transformed similarity  $f(S)$  will continue to admit an LSH if  $f(\cdot)$  is a probability generating function. The existence of an LSH that admits a similarity measure  $S$  will therefore give rise to the existence of solutions to the approximate similarity search problem for the much larger class of similarities  $f(S)$ . However, this still leaves open the problem of finding efficient explicit constructions, and as it turns out, the LSH property  $\Pr[h(\mathbf{x}) = h(\mathbf{y})] = S(\mathbf{x}, \mathbf{y})$ , while intuitively appealing and useful for similarity estimation,

does not necessarily imply that the LSH is optimal for solving the approximate search problem for the measure  $S$ . The problem of finding tight upper and lower bounds on the  $\rho$ -value that can be obtained through the LSH framework for data-independent  $(s_1, s_2)$ -approximate similarity search across the entire parameter space  $(s_1, s_2)$  remains open for two of the most common measures of set similarity: Jaccard similarity  $J(\mathbf{x}, \mathbf{y}) = |\mathbf{x} \cap \mathbf{y}|/|\mathbf{x} \cup \mathbf{y}|$  and cosine similarity  $C(\mathbf{x}, \mathbf{y}) = |\mathbf{x} \cap \mathbf{y}|/\sqrt{|\mathbf{x}||\mathbf{y}|}$ .

A random function from the MinHash family  $\mathcal{H}_{\text{minhash}}$  hashes a set  $\mathbf{x} \subset \{1, \dots, d\}$  to the first element of  $\mathbf{x}$  in a random permutation of the set  $\{1, \dots, d\}$ . For  $h \sim \mathcal{H}_{\text{minhash}}$  we have that  $\Pr[h(\mathbf{x}) = h(\mathbf{y})] = J(\mathbf{x}, \mathbf{y})$ , yielding an LSH solution to the approximate Jaccard similarity search problem. For cosine similarity the SimHash family  $\mathcal{H}_{\text{simhash}}$ , introduced by Charikar [13], works by sampling a random hyperplane in  $\mathbb{R}^d$  that passes through the origin and hashing  $\mathbf{x}$  according to what side of the hyperplane it lies on. For  $h \sim \mathcal{H}_{\text{simhash}}$  we have that  $\Pr[h(\mathbf{x}) = h(\mathbf{y})] = 1 - \arccos(C(\mathbf{x}, \mathbf{y}))/\pi$ , which can be used to derive a solution for cosine similarity, although not the clean solution that we could have hoped for in the style of MinHash for Jaccard similarity. There exists a number of different data-independent LSH approaches [2, 3, 34] that improve upon the  $\rho$ -value of SimHash. Perhaps surprisingly, it turns out that these approaches yield lower  $\rho$ -values for the  $(j_1, j_2)$ -approximate Jaccard similarity search problem compared to MinHash for certain combinations of  $(j_1, j_2)$ . Unfortunately, while asymptotically superior these techniques suffer from a non-trivial  $o_n(1)$ -term in the exponent that only decreases very slowly with  $n$ . In comparison, both MinHash and SimHash are simple to describe and have closed expressions for their  $\rho$ -values. Furthermore, MinHash and SimHash both have the advantage of being efficient in the sense that a hash function can be represented using space  $O(d)$  and the time to compute  $h(\mathbf{x})$  is  $O(|\mathbf{x}|)$ .

In Table 1 we show how the upper bounds for similarity search under different measures of set similarity relate to each other in the case where all sets are  $t$ -sparse. In addition to Hamming distance and Jaccard similarity, we consider Braun-Blanquet similarity [10] defined as

$$B(\mathbf{x}, \mathbf{y}) = |\mathbf{x} \cap \mathbf{y}|/\max(|\mathbf{x}|, |\mathbf{y}|), \quad (1)$$

which for  $t$ -sparse vectors is identical to cosine similarity. When the query and the sets in  $P$  can have different sizes the picture becomes muddled, and the question of which of the known algorithms is best for each measure  $S$  is complicated. In Section 5 we treat the problem of different set sizes and provide a brief discussion for Jaccard similarity, specifically in relation to our upper bound for Braun-Blanquet similarity.

Similarity search under set similarity and the batched version often referred to as *set similarity join* [7, 8] have also been studied extensively in the information retrieval and database literature, but mostly without providing theoretical guarantees on performance. Recently the notion of containment search, where the similarity measure is the (unnormalized) intersection size, was studied in the LSH framework [33]. This is a special case of *maximum inner product* search [1, 33]. However, these techniques do not give improvements in our setting.

**Similarity Estimation.** Finally, we mention that another application of MinHash [11, 12] is the (easier) problem of *similarity*

**Table 1: Overview of  $\rho$ -values for similarity search with Hamming vectors of equal weight  $t$ .**

Measure \ Ref.	Hamming $r_1 < r_2$	Braun-Blanquet $b_1 > b_2$	Jaccard $j_1 > j_2$
Bit-sampling LSH [23]	$r_1/r_2$	$\frac{1-b_1}{1-b_2}$	$\frac{1-j_1}{1+j_1} / \frac{1-j_2}{1+j_2}$
Minhash LSH [11]	$\log \frac{1-r_1}{1+r_1} / \log \frac{1-r_2}{1+r_2}$	$\log \frac{b_1}{2-b_1} / \log \frac{b_2}{2-b_2}$	$\log(j_1) / \log(j_2)$
Angular LSH [2]	$\frac{r_1}{r_2} \frac{1-r_2/2}{1-r_1/2}$	$\frac{1-b_1}{1+b_1} / \frac{1-b_2}{1+b_2}$	$\frac{1-j_1}{1+3j_1} / \frac{1-j_2}{1+3j_2}$
Data-dep. LSH [5]	$\frac{r_1}{r_2} \frac{1}{2-r_1/r_2}$	$\frac{1-b_1}{1+b_1-2b_2}$	$\frac{(1-j_1)(1+j_2)}{1-j_1j_2+3(j_1-j_2)}$
<b>Theorem 1.1</b>	$\log(1-r_1) / \log(1-r_2)$	$\log(b_1) / \log(b_2)$	$\log \frac{2j_1}{1+j_1} / \log \frac{2j_2}{1+j_2}$

**Notes:** While most results in the literature are stated for a single measure, the fixed weight restriction gives a 1-1 correspondence that makes it possible to express the results in terms of other similarity measures. Hamming distances are normalized by a factor  $2t$  to lie in  $[0; 1]$ . Lower order terms of  $\rho$ -values are suppressed, and for bit-sampling LSH we assume that the Hamming distances are small relative to the dimensionality of the space, i.e., that  $2r_1t/d = o(1)$ .

*estimation*, where the task is to condense each vector  $\mathbf{x}$  into a short signature  $s(\mathbf{x})$  in such a way that the similarity  $J(\mathbf{x}, \mathbf{y})$  can be estimated from  $s(\mathbf{x})$  and  $s(\mathbf{y})$ . A related similarity estimation technique was independently discovered by Cohen [17]. Thorup [35] has shown how to perform similarity estimation using just a small amount of randomness in the definition of the function  $s(\cdot)$ . In another direction, Mitzenmacher et al. [26] showed that it is possible to improve the performance of MinHash for similarity estimation when the Jaccard similarity is close to 1, but for smaller similarities it is known that succinct encodings of MinHash such as the one in [25] comes within a constant factor of the optimal space for storing  $s(\mathbf{x})$  [31]. Curiously, our improvement to MinHash in the context of similarity *search* comes when the similarity is neither too large nor too small. Our techniques do not seem to yield any improvement for the similarity *estimation* problem.

## 1.2 Contribution

We show the following upper bound for approximate similarity search under Braun-Blanquet similarity:

**THEOREM 1.1.** *For every choice of constants  $0 < b_2 < b_1 < 1$  we can solve the  $(b_1, b_2)$ -approximate similarity search problem under Braun-Blanquet similarity with query time  $O(|q|n^\rho \log n)$  and space usage  $O(n^{1+\rho} \log n + \sum_{\mathbf{x} \in P} |\mathbf{x}|)$  where  $\rho = \log(1/b_1) / \log(1/b_2)$ .*

In the case where the sets are  $t$ -sparse our Theorem 1.1 gives the first strict improvement on the  $\rho$ -value for approximate Jaccard similarity search compared to the data-independent LSH approaches of MinHash and Angular LSH. Figure 2 shows an example of the improvement for a slice of the parameter space. The improvement is based on a new locality-sensitive mapping that considers a specific random collection of length- $k$  paths on the vertex set  $\{1, \dots, d\}$ , and associates each vector  $\mathbf{x}$  with the paths in the collection that only visits vertices in  $\{i \mid x_i = 1\}$ . Our data structure exploits that similar vectors will be associated with a common path with constant probability, while vectors with low similarity have a negligible

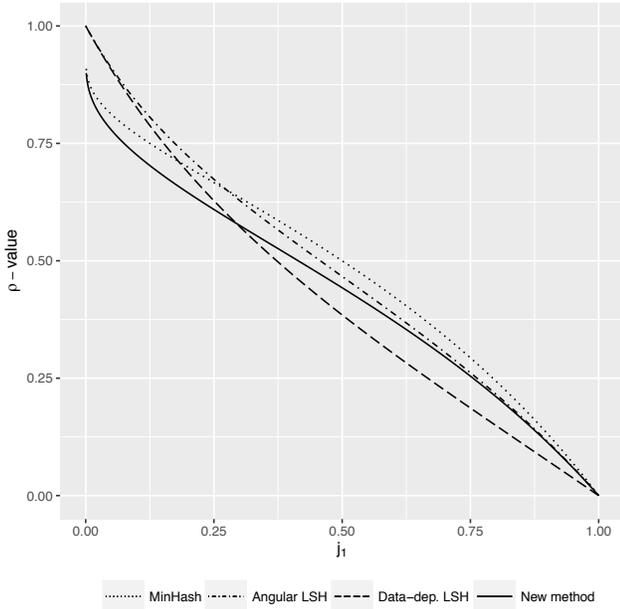
probability of sharing a path. However, the collection of paths has size superlinear in  $n$ , so an efficient method is required for locating the paths associated with a particular vector. Our choice of the collection of paths balances two opposing constraints: It is random enough to match the filtering performance of a truly random collection of sets, and at the same time it is structured enough to allow efficient search for sets matching a given vector. The search procedure is comparable in simplicity to the classical techniques of bit sampling, MinHash, SimHash, and  $p$ -stable LSH, and we believe it might be practical. This is in contrast to most theoretical advances in similarity search in the past ten years that suffer from  $o(1)$  terms in the exponent of complexity bounds.

**Intuition.** Recall that we will think of a vector  $\mathbf{x} \in \{0, 1\}^d$  also as a set,  $\{i \mid x_i = 1\}$ . MinHash can be thought of as a way of sampling an element  $i_{\mathbf{x}}$  from  $\mathbf{x}$ , namely, we let  $i_{\mathbf{x}} = \arg \min_{i \in \mathbf{x}} h(i)$  where  $h$  is a random hash function. For sets  $\mathbf{x}$  and  $\mathbf{y}$  the probability that  $i_{\mathbf{x}} = i_{\mathbf{y}}$  equals their Jaccard similarity  $J(\mathbf{x}, \mathbf{y})$ , which is much higher than if the samples had been picked independently. Consider the case in which  $|\mathbf{x}| = |\mathbf{y}| = t$ , so  $J(\mathbf{x}, \mathbf{y}) = \frac{|\mathbf{x} \cap \mathbf{y}|}{2t - |\mathbf{x} \cap \mathbf{y}|}$ . Another way of sampling is to compute  $I_{\mathbf{x}} = \mathbf{x} \cap \mathbf{b}$ , where  $\Pr[i \in \mathbf{b}] = 1/t$ , independently for each  $i \in [d]$ . The expected size of  $I_{\mathbf{x}}$  is 1, so this sample has the same expected ‘‘cost’’ as the MinHash-based sample. But if the Jaccard similarity is small, the latter samples are more likely to overlap:

$$\Pr[I_{\mathbf{x}} \cap I_{\mathbf{y}} \neq \emptyset] = 1 - (1 - 1/t)^{|\mathbf{x} \cap \mathbf{y}|} \approx 1 - e^{-|\mathbf{x} \cap \mathbf{y}|/t} \approx |\mathbf{x} \cap \mathbf{y}|/t,$$

nearly a factor of 2 improvement. In fact, whenever  $|\mathbf{x} \cap \mathbf{y}| < 0.6t$  we have  $\Pr[I_{\mathbf{x}} \cap I_{\mathbf{y}} \neq \emptyset] > \Pr[i_{\mathbf{x}} = i_{\mathbf{y}}]$ . So in a certain sense, MinHash is not the best way of collecting evidence for the similarity of two sets. (This observation is not new, and has been made before e.g. in [18].)

**Locality-Sensitive Maps.** The intersection of the samples  $I_{\mathbf{x}}$  does not correspond directly to hash collisions, so it is not clear how to turn this insight into an algorithm in the LSH framework. Instead, we will consider a generalization of both the locality sensitive



**Figure 2: Exponent when searching for a vector with Jaccard similarity  $j_1$  with approximation factor 2 (i.e., guaranteed to return a vector with Jaccard similarity  $j_1/2$ ) for various methods in the setting where all sets have the same size. Our new method is the best data-independent method, and is better than data-dependent LSH up to about  $j_1 \approx 0.3$ .**

filtering (LSF) and LSH frameworks where we define a distribution  $\mathcal{M}$  over maps  $M: \{0, 1\}^d \rightarrow 2^R$ . The map  $M$  performs the same task as the LSH data structure: It takes a vector  $\mathbf{x}$  and returns a set of memory locations  $M(\mathbf{x}) \subseteq \{1, \dots, R\}$ . A randomly sampled map  $M \sim \mathcal{M}$  has the property that if a pair of points  $\mathbf{x}, \mathbf{y}$  are close then  $M(\mathbf{x}) \cap M(\mathbf{y}) \neq \emptyset$  with constant probability, while if  $\mathbf{x}, \mathbf{y}$  are distant then the expected size  $|M(\mathbf{x}) \cap M(\mathbf{y})|$  is small (much smaller than 1). It is now straightforward to see that this distribution can be used to construct a data structure for similarity search by storing each data point  $\mathbf{x} \in P$  in the set of memory locations or buckets  $M(\mathbf{x})$ . A query for a point  $\mathbf{y}$  is performed by computing the similarity between  $\mathbf{y}$  and every point  $\mathbf{x}$  contained in the set buckets  $M(\mathbf{y})$ , reporting the first sufficiently similar point found.

**CHOSEN PATH.** It turns out that to most efficiently filter out vectors of low similarity in the setting where all sets have equal size, we would like to map each data point  $\mathbf{x} \in \{0, 1\}^d$  to a collection  $M(\mathbf{x})$  of random subsets of  $\{0, 1\}^d$  that are contained in  $\mathbf{x}$ . Furthermore, to best distinguish similar from dissimilar vectors when solving the approximate similarity search problem, we would like the random subsets of  $\{0, 1\}^d$  to have size  $\Theta(\log n)$ . This leads to another obstacle: The collection of subsets of  $\{0, 1\}^d$  required to ensure that  $M(\mathbf{x}) \cap M(\mathbf{y}) \neq \emptyset$  for similar points, i.e., that  $M$  maps to a subset contained in  $\mathbf{x} \cap \mathbf{y}$ , is very large. The space usage and evaluation time of a locality-sensitive map  $M$  to fully random subsets of  $\{0, 1\}^d$  would far exceed  $n$ , rendering the solution useless.

To overcome this we create the samples in a gradual, correlated way using a pairwise independent branching process that turns out to yield “sufficiently random” samples for the argument to go through.

**Lower Bound.** On the lower bound side we show that our solution for Braun-Blanquet similarity is best possible in terms of parameters  $b_1$  and  $b_2$  within the class of solutions that can be characterized as data-independent locality-sensitive maps. The lower bound works by showing that a family of locality-sensitive maps for Braun-Blanquet similarity with a  $\rho$ -value below  $\log(1/b_1)/\log(1/b_2)$  can be used to construct a locality-sensitive hash family for the  $c$ -approximate near neighbor problem in Hamming space with a  $\rho$ -value below  $1/c$ , thereby contradicting the LSH lower bound by O’Donnell et al. [30]. We state the lower bound here in terms of locality-sensitive hashing, formally defined in Section 2.

**THEOREM 1.2.** *For every choice of constants  $0 < b_2 < b_1 < 1$  any  $(b_1, b_2, p_1, p_2)$ -sensitive hash family  $\mathcal{H}_B$  for  $\{0, 1\}^d$  under Braun-Blanquet similarity must satisfy*

$$\rho(\mathcal{H}_B) = \frac{\log(1/p_1)}{\log(1/p_2)} \geq \frac{\log(1/b_1)}{\log(1/b_2)} - O\left(\frac{\log(d/p_2)}{d}\right)^{1/3}.$$

The details showing how this LSH lower bound implies a lower bound for locality-sensitive maps are given in Section 4.

## 2 PRELIMINARIES

As stated above we will view  $\mathbf{x} \in \{0, 1\}^d$  both as a vector and as a subset of  $[d] = \{1, \dots, d\}$ . Define  $\mathbf{x}$  to be  $t$ -sparse if  $|\mathbf{x}| = t$ ; we will be interested in the setting where  $t \leq d/2$ , and typically the sparse setting  $t \ll d$ . Although many of the concepts we use hold for general spaces, for simplicity we state definitions in the same setting as our results: the boolean hypercube  $\{0, 1\}^d$  under some measure of similarity  $S: \{0, 1\}^d \times \{0, 1\}^d \rightarrow [0, 1]$ .

**Definition 2.1.** (Approximate similarity search) Let  $P \subset \{0, 1\}^d$  be a set of  $|P| = n$  data vectors, let  $S: \{0, 1\}^d \times \{0, 1\}^d \rightarrow [0, 1]$  be a similarity measure, and let  $s_1, s_2 \in [0, 1]$  such that  $s_1 > s_2$ . A solution to the  $(s_1, s_2)$ - $S$ -similarity search problem is a data structure that supports the following query operation: on input  $\mathbf{q} \in \{0, 1\}^d$  for which there exists a vector  $\mathbf{x} \in P$  with  $S(\mathbf{x}, \mathbf{q}) \geq s_1$ , return  $\mathbf{x}' \in P$  with  $S(\mathbf{x}', \mathbf{q}) > s_2$ .

Our data structures are randomized, and queries succeed with probability at least  $1/2$  (the probability can be made arbitrarily close to 1 by independent repetition). Sometimes similarity search is formulated as searching for vectors that are near  $\mathbf{q}$  according to the distance measure  $D(\mathbf{x}, \mathbf{y}) = 1 - S(\mathbf{x}, \mathbf{y})$ . For our purposes it is natural to phrase conditions in terms of similarity, but we compare to solutions originally described as “near neighbor” methods.

Many of the best known solutions to approximate similarity search problems are based on a technique of randomized space partitioning. This technique has been formalized in the locality-sensitive hashing framework [23] and the closely related locality-sensitive filtering framework [9, 16].

**Definition 2.2.** (Locality-sensitive hashing [23]) A  $(s_1, s_2, p_1, p_2)$ -sensitive family of hash functions for a similarity measure  $S$  on  $\{0, 1\}^d$  is a distribution  $\mathcal{H}_S$  over functions  $h: \{0, 1\}^d \rightarrow R$  such

that for all  $\mathbf{x}, \mathbf{y} \in \{0, 1\}^d$  and random  $h$  sampled according to  $\mathcal{H}_S$ : If  $S(\mathbf{x}, \mathbf{y}) \geq s_1$  then  $\Pr[h(\mathbf{x}) = h(\mathbf{y})] \geq p_1$ , and if  $S(\mathbf{x}, \mathbf{y}) \leq s_2$  then  $\Pr[h(\mathbf{x}) = h(\mathbf{y})] \leq p_2$ .

The range  $R$  of the family will typically be fairly small such that an element of  $R$  can be represented in a constant number of machine words. In the following we assume for simplicity that the family of hash functions is *efficient* such that  $h(\mathbf{x})$  can be computed in time  $O(|\mathbf{x}|)$ . Furthermore, we will assume that the time to compute the similarity  $S(\mathbf{x}, \mathbf{y})$  can be upper bounded by the time it takes to compute the size of the intersection of preprocessed sets, i.e.,  $O(\min(|\mathbf{x}|, |\mathbf{y}|))$ .

Given a locality-sensitive family it is quite simple to obtain a solution to the approximate similarity search problem, essentially by hashing points to buckets such that close points end up in the same bucket while distant points are kept apart.

**LEMMA 2.3 (LSH FRAMEWORK [21, 23]).** *Given a  $(s_1, s_2, p_1, p_2)$ -sensitive family of hash functions it is possible to solve the  $(s_1, s_2)$ - $S$ -similarity search problem with query time  $O(|\mathbf{q}|n^\rho \log n)$  and space usage  $O(n^{1+\rho} + \sum_{\mathbf{x} \in P} |\mathbf{x}|)$  where  $\rho = \log(1/p_1)/\log(1/p_2)$ .*

The upper bound presented in this paper does not quite fit into the existing frameworks. However, we would like to apply existing LSH lower bound techniques to our algorithm. Therefore we define a more general framework that captures solutions constructed using the LSH and LSF framework, as well as the upper bound presented in this paper.

**Definition 2.4 (Locality-sensitive map).** A  $(s_1, s_2, m_1, m_2)$ -sensitive family of maps for a similarity measure  $S$  on  $\{0, 1\}^d$  is a distribution  $\mathcal{M}_S$  over mappings  $M: \{0, 1\}^d \rightarrow 2^R$  (where  $2^R$  denotes the power set of  $R$ ) such that for all  $\mathbf{x}, \mathbf{y} \in \{0, 1\}^d$  and random  $M \in \mathcal{M}_S$ :

- (1)  $\mathbb{E}[|M(\mathbf{x})|] \leq m_1$ .
- (2) If  $S(\mathbf{x}, \mathbf{y}) \leq s_2$  then  $\mathbb{E}[|M(\mathbf{x}) \cap M(\mathbf{y})|] \leq m_2$ .
- (3) If  $S(\mathbf{x}, \mathbf{y}) \geq s_1$  then  $\Pr[M(\mathbf{x}) \cap M(\mathbf{y}) \neq \emptyset] \geq 1/2$ .

Once we have a family of locality-sensitive maps  $\mathcal{M}$  we can use it to obtain a solution to the  $(s_1, s_2)$ - $S$ -similarity search problem.

**LEMMA 2.5.** *Given a  $(s_1, s_2, m_1, m_2)$ -sensitive family of maps  $\mathcal{M}$  we can solve the  $(s_1, s_2)$ - $S$ -similarity search problem with query time  $O(m_1 + nm_2|\mathbf{q}| + T_M)$  and space usage  $O(nm_1 + \sum_{\mathbf{x} \in P} |\mathbf{x}|)$  where  $T_M$  is the time to evaluate a map  $M \in \mathcal{M}$ .*

**PROOF.** We construct the data structure by sampling a map  $M$  from  $\mathcal{M}$  and use it to place points in  $P$  into buckets. To run a query for a point  $\mathbf{q}$  we proceed by evaluating  $M(\mathbf{q})$  and computing the similarity between  $\mathbf{q}$  and the points in the buckets associated with  $M(\mathbf{q})$ . If a sufficiently similar point is found we return it. We get rid of the expectation in the guarantees by independent repetitions and applying Markov's inequality.  $\square$

**Model of Computation.** We assume the standard word RAM model [20] with word size  $\Theta(\log n)$ , where  $n = |P|$ . In order to be able to draw random functions from a family of functions we augment the model with an instruction that generates a machine word uniformly at random in constant time.

### 3 UPPER BOUND

We will describe a family of locality-sensitive maps  $\mathcal{M}_B$  for solving the  $(b_1, b_2)$ - $B$ -similarity search problem, where  $B$  is Braun-Blanquet similarity (1). After describing  $\mathcal{M}_B$  we will give an efficient implementation of  $M \in \mathcal{M}_B$  and show how to set parameters to obtain our Theorem 1.1.

#### 3.1 Chosen Path

The CHOSEN PATH family  $\mathcal{M}_B$  is defined by  $k$  random hash functions  $h_1, \dots, h_k$  where  $h_i: [w] \times [d]^i \rightarrow [0, 1]$  and  $w$  is a positive integer. The evaluation of a map  $M_k \in \mathcal{M}_B$  proceeds in a sequence of  $k + 1$  steps that can be analyzed as a Galton-Watson branching process, originally devised to investigate population growth under the assumption of identical and independent offspring distributions. In the first step  $i = 0$  we create a population of  $w$  starting points

$$M_0(\mathbf{x}) = [w]. \quad (2)$$

In subsequent steps, every path that has survived so far produces offspring according to a random process that depends on  $h_i$  and the element  $\mathbf{x} \in \{0, 1\}^d$  being evaluated. We use  $p \circ j$  to denote concatenation of a path  $p$  with a vertex  $j$ .

$$M_i(\mathbf{x}) = \left\{ p \circ j \mid p \in M_{i-1}(\mathbf{x}) \wedge h_i(p \circ j) < \frac{x_j}{b_1|\mathbf{x}|} \right\}. \quad (3)$$

Observe that  $h_i(p \circ j) < \frac{x_j}{b_1|\mathbf{x}|}$  can only hold when  $x_j = 1$ , so the paths in  $M_i(\mathbf{x})$  are constrained to  $j \in \mathbf{x}$ . The set  $M(\mathbf{x}) = M_k(\mathbf{x})$  is given by the paths that survive to the  $k$ th step. We will proceed by bounding the evaluation time of  $M \in \mathcal{M}_B$  as well as showing the locality-sensitive properties of  $\mathcal{M}_B$ . In particular, for similar points  $\mathbf{x}, \mathbf{y} \in \{0, 1\}^d$  with  $B(\mathbf{x}, \mathbf{y}) \geq b_1$  we will show that with probability at least  $1/2$  there will be a path that is chosen by both  $\mathbf{x}$  and  $\mathbf{y}$ .

**LEMMA 3.1 (PROPERTIES OF CHOSEN PATH).** *For all  $\mathbf{x}, \mathbf{y} \in \{0, 1\}^d$ , integer  $i \geq 0$ , and random  $M \in \mathcal{M}_B$ :*

- (1)  $\mathbb{E}[|M_i(\mathbf{x})|] \leq (1/b_1)^i w$ .
- (2) If  $B(\mathbf{x}, \mathbf{y}) < b_2$  then  $\mathbb{E}[|M_i(\mathbf{x}) \cap M_i(\mathbf{y})|] \leq (b_2/b_1)^i w$ .
- (3) If  $B(\mathbf{x}, \mathbf{y}) \geq b_1$  then  $\Pr[M_i(\mathbf{x}) \cap M_i(\mathbf{y}) \neq \emptyset] \geq i/(i + w)$ .

**PROOF.** We prove each property by induction on  $i$ . The base cases  $i = 0$  follow from (2). Now consider the inductive step for property 1. Let  $\mathbb{1}\{\mathcal{P}\}$  denote the indicator function for predicate  $\mathcal{P}$ . Using independence of the hash functions  $h_i$  we get:

$$\begin{aligned} \mathbb{E}[|M_i(\mathbf{x})|] &= \mathbb{E} \left[ \sum_{p \in M_{i-1}(\mathbf{x})} \sum_{j \in [d]} \mathbb{1} \left\{ h_i(p \circ j) < \frac{x_j}{b_1|\mathbf{x}|} \right\} \right] \\ &= \mathbb{E} \left[ \sum_{p \in M_{i-1}(\mathbf{x})} 1 \right] \mathbb{E} \left[ \sum_{j \in [d]} \mathbb{1} \left\{ h_i(p \circ j) < \frac{x_j}{b_1|\mathbf{x}|} \right\} \right] \\ &\leq \mathbb{E}[|M_{i-1}(\mathbf{x})|]/b_1 \\ &\leq (1/b_1)^i w. \end{aligned}$$

The last inequality uses the induction hypothesis. We use the same approach for the second property where we let  $X_i = M_i(\mathbf{x}) \cap M_i(\mathbf{y})$ .

$$\begin{aligned} \mathbb{E}[|X_i|] &= \mathbb{E} \left[ \sum_{p \in X_{i-1}} \sum_{j \in [d]} \mathbb{1} \left\{ h_i(p \circ j) < \frac{x_j}{b_1 |\mathbf{x}|} \wedge h_i(p \circ j) < \frac{y_j}{b_1 |\mathbf{y}|} \right\} \right] \\ &= \mathbb{E} \left[ \sum_{p \in X_{i-1}} 1 \right] \sum_{j \in [d]} \Pr \left[ h_i(p \circ j) < \frac{\min(x_j, y_j)}{b_1 \max(|\mathbf{x}|, |\mathbf{y}|)} \right] \\ &\leq \mathbb{E}[|X_{i-1}|] (B(\mathbf{x}, \mathbf{y})/b_1) \\ &\leq (B(\mathbf{x}, \mathbf{y})/b_1)^i w . \end{aligned}$$

To prove the third property we bound the variance of  $|X_i|$  and apply Chebyshev's inequality to bound the probability of  $X_i = \emptyset$ . First consider the case where  $|\mathbf{x}| \leq 1/b_1$  and  $|\mathbf{y}| \leq 1/b_1$ . Here it must hold that  $X_i > 0$  as intersecting paths exist ( $b_1 > 0$ ) and always activate. In all other cases we have that

$$\mathbb{E}[|X_i|] = (B(\mathbf{x}, \mathbf{y})/b_1)^i w .$$

Knowing the expected value we can apply Chebyshev's inequality once we have an upper bound for  $\text{Var}[|X_i|] = \mathbb{E}[|X_i|^2] - \mathbb{E}[|X_i|]^2$ . Specifically we show that  $\mathbb{E}[|X_i|^2] \leq wi(B(\mathbf{x}, \mathbf{y})/b_1)^{2i}$ , by induction on  $i$ . To simplify notation we define the indicator variable

$$Y_{p,j} = \mathbb{1} \left\{ h_i(p \circ j) < \frac{x_j}{b_1 |\mathbf{x}|} \wedge h_i(p \circ j) < \frac{y_j}{b_1 |\mathbf{y}|} \right\}$$

where we suppress the subscript  $i$ . First observe that

$$\mathbb{E}[Y_{p,j}] = 1/(b_1 \max(|\mathbf{x}|, |\mathbf{y}|)) .$$

By (3) we see that  $|X_i| = \sum_{p \in X_{i-1}} \sum_{j \in [d]} Y_{p,j}$ , which means:

$$\begin{aligned} \mathbb{E}[|X_i|^2] &= \mathbb{E} \left[ \left( \sum_{p \in X_{i-1}} \sum_{j \in [d]} Y_{p,j} \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{p \in X_{i-1}} \sum_{j \in [d]} Y_{p,j}^2 \right] \\ &\quad + \mathbb{E} \left[ \sum_{p, p' \in X_{i-1}} \sum_{j, j' \in [d]} Y_{p,j} Y_{p',j'} \mathbb{1} \{ (p, j) \neq (p', j') \} \right] \\ &< \mathbb{E}[|X_{i-1}|] (B(\mathbf{x}, \mathbf{y})/b_1) + \mathbb{E}[|X_{i-1}|^2] (B(\mathbf{x}, \mathbf{y})/b_1)^2 \\ &\leq \sum_{s=1}^i \mathbb{E}[|X_{i-s}|] (B(\mathbf{x}, \mathbf{y})/b_1)^{2s-1} + \mathbb{E}[|X_0|^2] (B(\mathbf{x}, \mathbf{y})/b_1)^{2i} \\ &= \mathbb{E}[|X_i|] \sum_{s=0}^{i-1} (B(\mathbf{x}, \mathbf{y})/b_1)^s + \mathbb{E}[|X_i|]^2 \\ &\leq wi(B(\mathbf{x}, \mathbf{y})/b_1)^{2i} + \mathbb{E}[|X_i|]^2 . \end{aligned}$$

The third property now follows from a one-sided version of Chebyshev's inequality applied to  $|X_i|$ .  $\square$

### 3.2 Implementation Details

Lemma 3.1 continues to hold when the hash functions  $h_1, \dots, h_k$  are individually 2-independent (and mutually independent) since we only use bounds on the first and second moment of the hash values. We can therefore use a simple and practical scheme such as Zobrist hashing [37] that hashes strings of  $\Theta(\log n)$  bits to strings of

$\Theta(\log n)$  bits in  $O(1)$  time using space, say,  $O(n^{1/2})$ . It is not hard to see that the domain and range of  $h_1, \dots, h_k$  can be compressed to  $O(\log n)$  bits (causing a negligible increase in the failure probability of the data structure). We simply hash the paths  $p \in M_i(\mathbf{x})$  to intermediate values of  $O(\log n)$  bits, avoiding collisions with high probability, and in a similar vein, with high probability  $O(\log n)$  bits of precision suffice to determine whether  $h_i(p \circ j) < \frac{x_j}{b_1 |\mathbf{x}|}$ .

We now consider how to parameterize  $\mathcal{M}_B$  to solve the  $(b_1, b_2)$ - $B$ -similarity problem on a set  $P$  of  $|P| = n$  points for every choice of constant parameters  $0 < b_2 < b_1 < 1$ , independent of  $n$ . Note that we exclude  $b_1 = 1$  (which would correspond to identical vectors that can be found in time  $O(1)$  by resorting to standard hashing) and  $b_2 = 0$  (for which every data point would be a valid answer to a query). We set parameters

$$\begin{aligned} k &= \lceil \log(n)/\log(1/b_2) \rceil, \\ w &= 2k \end{aligned}$$

from which it follows that  $\mathcal{M}_B$  is  $(b_1, b_2, m_1, m_2)$ -sensitive with  $m_1 = n^\rho w/b_1$  and  $m_2 = n^{\rho-1} w$  where  $\rho = \log(1/b_1)/\log(1/b_2)$ . To bound the expected evaluation time of  $M_k$  we use Zobrist hashing as well as intermediate hashes for the paths as described above. In the  $i$ th step in the branching process the expected number of hash function evaluations is bounded by  $|\mathbf{q}|$  times the number of paths alive at step  $i-1$ . We can therefore bound the expected time to compute  $M_k(\mathbf{q})$  by

$$\sum_{i=0}^{k-1} \mathbb{E}[|\mathbf{q}| |M_i(\mathbf{q})|] \leq \frac{b_1^{-k} - 1}{b_1^{-1} - 1} |\mathbf{q}| w = O(|\mathbf{q}| n^\rho w). \quad (4)$$

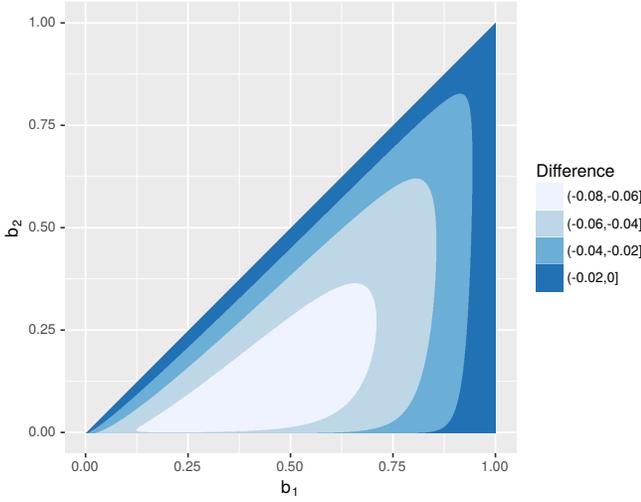
This completes the proof of Theorem 1.1.<sup>1</sup>

### 3.3 Comparison

We will proceed by comparing our Theorem 1.1 to results that can be achieved using existing techniques. Again we focus on the setting where data points and query points are exactly  $t$ -sparse. An overview of different techniques for three measures of similarity is shown in Table 1. To summarize: The CHOSEN PATH algorithm of Theorem 1.1 improves upon all existing data-independent results over the entire  $0 < b_2 < b_1 < 1$  parameter space. Furthermore, we improve upon the best known *data-dependent* techniques [5] for a large part of the parameter space (see Figure 5). The details of the comparisons are given in Appendix B.

**MinHash.** For  $t$ -sparse vectors there is a 1-1 mapping between Braun-Blanquet and Jaccard similarity. In this setting  $J(\mathbf{x}, \mathbf{y}) = B(\mathbf{x}, \mathbf{y})/(2 - B(\mathbf{x}, \mathbf{y}))$ . Let  $b_1 = 2j_1/(j_1 + 1)$  and  $b_2 = 2j_2/(j_2 + 1)$  be the Braun-Blanquet similarities corresponding to Jaccard similarities  $j_1$  and  $j_2$ . The LSH framework using MinHash achieves  $\rho_{\text{minhash}} = \log\left(\frac{b_1}{2-b_1}\right) / \log\left(\frac{b_2}{2-b_2}\right)$ ; this should be compared to  $\rho = \log(b_1)/\log(b_2)$  achieved in Theorem 1.1. Since the function  $f(z) = \log(\frac{z}{2-z})/\log z$  is monotonically increasing in  $[0, 1]$  we have that  $\rho/\rho_{\text{minhash}} = f(b_2)/f(b_1) < 1$ , i.e.,  $\rho$  is always smaller than  $\rho_{\text{minhash}}$ . As an example, for  $j_1 = 0.2$  and  $j_2 = 0.1$  we get

<sup>1</sup>We know of a way of replacing the multiplicative factor  $|\mathbf{q}|$  in equation (4) by an additive term of  $O(|\mathbf{q}|k)$  by choosing the hash functions  $h_i$  carefully, but do not discuss this improvement here since  $|\mathbf{q}|$  can be assumed to be polylogarithmic and our focus is on the exponent of  $n$ .



**Figure 3: The difference  $\rho - \rho_{\text{minhash}}$  comparing CHOSEN PATH and MinHash in terms of Braun-Blanquet similarities  $0 < b_2 < b_1 < 1$ .**

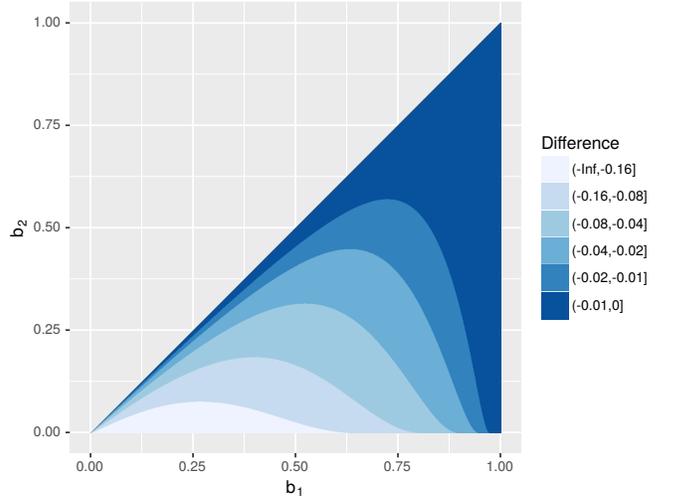
$\rho = 0.644\dots$  while  $\rho_{\text{minhash}} = 0.698\dots$  Figure 3 shows the difference for the whole parameter space.

**Angular LSH.** Since our vectors are exactly  $t$ -sparse Braun-Blanquet similarities correspond directly to dot products (which in turn correspond to angles). Thus we can apply angular LSH such as SimHash [13] or cross-polytope LSH [2]. As observed in [16] one can express the  $\rho$ -value of cross-polytope LSH in terms of dot products as  $\rho_{\text{angular}} = \frac{1-b_1}{1+b_1} / \frac{1-b_2}{1+b_2}$ . Since the function  $f'(z) = (1+z)\log(z)/(1-z)$  is negative and monotonically increasing in  $[0; 1]$  we have that  $\rho/\rho_{\text{angular}} = f'(b_1)/f'(b_2) < 1$ , i.e.,  $\rho$  is always smaller than  $\rho_{\text{angular}}$ . In the above example, for  $j_1 = 0.2$  and  $j_2 = 0.1$  we have  $\rho_{\text{angular}} = 0.722\dots$  which is about 0.078 more than CHOSEN PATH. See Figure 4 for a visualization of the difference for the whole parameter space.

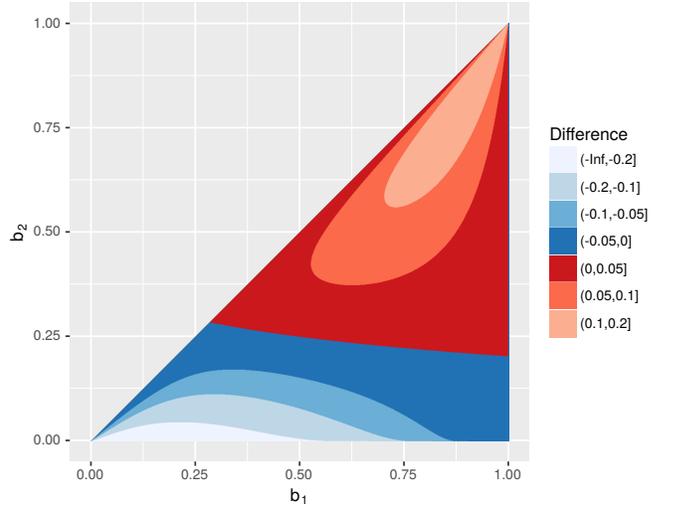
**Data-Dependent Hamming LSH.** The Hamming distance between two  $t$ -sparse vectors with Braun-Blanquet similarity  $b$  is  $2t(1-b)$ , since the intersection of the vectors has size  $tb$ . This means that  $(b_1, b_2)$ - $B$ -similarity search can be reduced to Hamming similarity search with approximation factor  $c = (2t(1-b_1))/(2t(1-b_2)) = (1-b_1)/(1-b_2)$ . As mentioned above, the *data dependent* LSH technique of [5] achieves  $\rho = 1/(2c-1)$  ignoring  $o_n(1)$  terms. In terms of  $b_1$  and  $b_2$  this is  $\rho_{\text{datadep}} = \frac{1-b_1}{1+b_1-2b_2}$ , which is incomparable to the  $\rho$  of Theorem 1.1. In Appendix B we show that  $\rho < \rho_{\text{datadep}}$  whenever  $b_2 \leq 1/5$ , or equivalently, whenever  $j_2 \leq 1/9$ . Revisiting the above example, for  $j_1 = 0.2$  and  $j_2 = 0.1$  we have  $\rho_{\text{datadep}} = 0.6875$  which is about 0.043 more than CHOSEN PATH. Figure 5 gives a comparison covering the whole parameter space.

## 4 LOWER BOUND

In this section we will show a locality-sensitive hashing lower bound for  $\{0, 1\}^d$  under Braun-Blanquet similarity. We will first show that LSH lower bounds apply to the class of solutions to the



**Figure 4: The difference  $\rho - \rho_{\text{angular}}$  comparing CHOSEN PATH and angular LSH in terms of Braun-Blanquet similarities  $0 < b_2 < b_1 < 1$ .**



**Figure 5: The difference  $\rho - \rho_{\text{datadep}}$  comparing CHOSEN PATH and data-dependent LSH in terms of Braun-Blanquet similarities  $0 < b_2 < b_1 < 1$ . In the area of the parameter space that is colored blue we have that  $\rho \leq \rho_{\text{datadep}}$  while for the red area it holds that  $\rho > \rho_{\text{datadep}}$ .**

approximate similarity search problem that are based on locality-sensitive maps, thereby including our own upper bound. Next we will introduce some relevant tools from the literature, in particular the LSH lower bounds for Hamming space by O'Donnell et al. [30] which we use, through a reduction, to show LSH lower bounds under Braun-Blanquet similarity.

**Lower Bounds For Locality-Sensitive Maps.** Because our upper bound is based on a locality-sensitive map  $\mathcal{M}_B$  and not

LSH-based we first show that LSH lower bounds apply to LSM-based solutions. This is not too surprising as both the LSH and LSF frameworks produce LSM-based solutions. We note that the idea of showing lower bounds for a more general class of algorithms that encompasses both LSH and LSF was used by Andoni et al. [4] in their list-of-points data structure lower bound for the space-time tradeoff of solutions to the approximate near neighbor problem in the random data regime. We use the approach of Christiani [16] to convert an LSM family into an LSH family using MinHash.

**LEMMA 4.1.** *Suppose we have a  $(s_1, s_2, m_1, m_2)$ -sensitive family of maps  $\mathcal{M}$  for a similarity measure  $S$  on  $\{0, 1\}^d$ . Then we can construct a  $(s_1, s_2, p_1, p_2)$ -sensitive family of hash functions  $\mathcal{H}$  for  $S$  such that  $p_1 = 1/8m$  and  $p_2 = m_2/m$  where  $m = \lceil 8m_1 \rceil$ .*

**PROOF.** We sample a function  $h$  from  $\mathcal{H}$  by sampling a function  $M$  from  $\mathcal{M}$ , modify  $M$  to output a set of fixed size, and apply MinHash to the resulting set. For  $M \in \mathcal{M}$  we define the function  $\tilde{M}$  where we ensure that the size of the output set is  $m$ . We note that the purpose of this step is to be able to simultaneously lower bound  $p_1$  and upper bound  $p_2$  for  $\mathcal{H}$  when we apply MinHash to the resulting sets.

$$\tilde{M}(\mathbf{x}) = \begin{cases} \{(\mathbf{x}, 1), \dots, (\mathbf{x}, m)\} & \text{if } |M(\mathbf{x})| \geq m, \\ \{(\mathbf{x}, 1), \dots, (\mathbf{x}, m - |M(\mathbf{x})|)\} \cup M(\mathbf{x}) & \text{otherwise.} \end{cases}$$

We proceed by applying MinHash to the set  $\tilde{M}(\mathbf{x})$ . Let  $\pi$  denote a random permutation of the range of  $\tilde{M}$  and define

$$h(\mathbf{x}) = \arg \min_{z \in \tilde{M}(\mathbf{x})} \pi(z).$$

We then have

$$\Pr[h(\mathbf{x}) = h(\mathbf{y})] = \sum_{\xi} \Pr[J(\tilde{M}(\mathbf{x}), \tilde{M}(\mathbf{y})) = \xi] \cdot \xi$$

summing over the finite set of all possible Jaccard similarities  $\xi = a/b$  with  $a, b \in \{0, 1, \dots, 2m\}$ . It is now fairly simple to lower bound  $p_1$  and upper bound  $p_2$ . Assume that  $\mathbf{x}, \mathbf{y}$  satisfy that  $S(\mathbf{x}, \mathbf{y}) \geq s_1$ . To lower bound  $p_1$  we use a union bound together with Markov's inequality to bound the following probability:

$$\begin{aligned} \Pr[\tilde{M}(\mathbf{x}) \cap \tilde{M}(\mathbf{y}) = \emptyset] & \leq \Pr[M(\mathbf{x}) \cap M(\mathbf{y}) = \emptyset \wedge |M(\mathbf{x})| \geq m \wedge |M(\mathbf{y})| \geq m] \\ & \leq \Pr[M(\mathbf{x}) \cap M(\mathbf{y}) = \emptyset] + \Pr[|M(\mathbf{x})| \geq m] + \Pr[|M(\mathbf{y})| \geq m] \\ & \leq 1/2 + 1/8 + 1/8 \end{aligned}$$

We therefore have that  $\Pr[\tilde{M}(\mathbf{x}) \cap \tilde{M}(\mathbf{y}) \neq \emptyset] \geq 1/4$ . In the event of a nonempty intersection the probability of collision is given by  $J(\tilde{M}(\mathbf{x}) \cap \tilde{M}(\mathbf{y})) \geq 1/2m$  allowing us to conclude that  $p_1 \geq 1/8m$ .

Bounding the collision probability for distant pairs of points  $\mathbf{x}, \mathbf{y}$  with  $S(\mathbf{x}, \mathbf{y}) \leq s_2$  we get

$$\sum_{\xi} \Pr[J(\tilde{M}(\mathbf{x}), \tilde{M}(\mathbf{y})) = \xi] \cdot \xi \leq (1/m) \sum_{i=1}^{\infty} \Pr[|\tilde{M}(\mathbf{x}) \cap \tilde{M}(\mathbf{y})| = i] = \frac{m_2}{m}.$$

□

We are now ready to justify the statement that LSH lower bounds apply to LSM, allowing us to restrict our attention to proving LSH lower bounds for Braun-Blanquet similarity.

**COROLLARY 4.2.** *Suppose that we have an LSM-based solution to the  $(s_1, s_2)$ - $S$ -similarity search problem with query time  $O(n^\rho)$ . Then there exists a family  $\mathcal{H}$  of locality-sensitive hash functions with  $\rho(\mathcal{H}) = \rho + O(1/\log n)$ .*

**PROOF.** The existence of the LSM-based solution implies that for every  $n$  there exists a  $(s_1, s_2, m_1, m_2)$ -sensitive family of maps  $\mathcal{M}$  with  $m_1 = O(n^\rho)$  and  $nm_2 = O(n^\rho)$ . The upper bound on  $\rho$  follows from applying Lemma 4.1. □

**LSH Lower Bounds for Hamming Space.** There exist a number of powerful results that lower bound the  $\rho$ -value that is attainable by locality-sensitive hashing and related approaches in various settings [4, 6, 16, 28, 30, 32]. O'Donnell et al. [30] showed an LSH lower bound of  $\rho = \log(1/p_1)/\log(1/p_2) \geq 1/c - o_d(1)$  for  $d$ -dimensional Hamming space under the assumption that  $p_2$  is not too small compared to  $d$ , i.e.,  $\log(1/p_2) = o(d)$ . The lower bound by O'Donnell et al. holds for  $(r, cr, p_1, p_2)$ -sensitive families for a particular choice of  $r$  that depends on  $d, p_2$ , and  $c$ , and where  $r$  is small compared to  $d$  (for instance, we have that  $r = \tilde{\Theta}(d^{2/3})$  when  $c$  and  $p_2$  are constant).

We state a simplified version of the lower bound due to O'Donnell et al. where  $r = \sqrt{d}$  that we will use as a tool to prove our lower bound for Braun-Blanquet similarity. The full proof of Lemma 4.3 is given in Appendix A.

**LEMMA 4.3.** *For every  $d \in \mathbb{N}$ ,  $1/d \leq p_2 \leq 1 - 1/d$ , and  $1 \leq c \leq d^{1/8}$  every  $(\sqrt{d}, c\sqrt{d}, p_1, p_2)$ -sensitive hash family  $\mathcal{H}$  for  $\{0, 1\}^d$  under Hamming distance must have*

$$\rho(\mathcal{H}) = \frac{\log(1/p_1)}{\log(1/p_2)} \geq \frac{1}{c} - O(d^{-1/4}). \quad (5)$$

In general, good lower bounds for the entire parameter space  $(r, cr)$  are not known, although the techniques by O'Donnell et al. appear to yield a bound of  $\rho \gtrsim \log(1 - 2r/d)/\log(1 - 2cr/d)$ . This is far from tight as can be seen by comparing it to the bit-sampling [23] upper bound of  $\rho = \log(1 - r/d)/\log(1 - cr/d)$ . Existing lower bounds are tight in two different settings. First, in the setting where  $cr \approx d/2$  (random data), lower bounds [6, 19, 28] match various instantiations of angular LSH [2, 3, 34]. Second, in the setting where  $r \ll d$ , the lower bound by O'Donnell et al. [30] becomes  $\rho \gtrsim \log(1 - 2r/d)/\log(1 - 2cr/d) \approx 1/c$ , matching bit-sampling LSH [23] as well as Angular LSH.

## 4.1 Braun-Blanquet LSH Lower Bound

We are now ready to prove the LSH lower bound from Theorem 1.2. The lower bound together with Corollary 4.2 shows that the  $\rho$ -value of Theorem 1.1 is best possible up to  $o_d(1)$  terms within the class of data-independent locality-sensitive maps for Braun-Blanquet similarity. Furthermore, the lower bound also applies to angular distance on the unit sphere where it comes close to matching the best known upper bounds for much of the parameter space as can be seen from Figure 4.

**Proof Sketch.** The proof works by assuming the existence of a  $(b_1, b_2, p_1, p_2)$ -sensitive family  $\mathcal{H}_B$  for  $\{0, 1\}^d$  under Braun-Blanquet similarity with  $\rho = \log(1/b_1)/\log(1/b_2) - \gamma$  for some  $\gamma > 0$ . We use a transformation  $T$  from Hamming space to Braun-Blanquet

similarity to show that the existence of  $\mathcal{H}_B$  implies the existence of a  $(r, cr, p'_1, p'_2)$ -sensitive family  $\mathcal{H}_H$  for  $D$ -dimensional Hamming space that will contradict the lower bound of O'Donnell et al. [30] as stated in Lemma 4.3 for some appropriate choice of  $\gamma = \gamma(d, p_2)$ .

We proceed by giving an informal description of a simple ‘‘tensoring’’ technique for converting a similarity search problem in Hamming space into a Braun-Blanquet set similarity problem for target similarity thresholds  $b_1, b_2$ . For  $\mathbf{x} \in \{0, 1\}^d$  define

$$\tilde{\mathbf{x}} = \{(i, \mathbf{x}_i) \mid i \in [d]\}$$

and for a positive integer  $\tau$  define  $\mathbf{x}^{\otimes \tau} = \{(v_1, \dots, v_\tau) \mid v_i \in \tilde{\mathbf{x}}\}$ . We have that  $|\mathbf{x}^{\otimes \tau}| = |\tilde{\mathbf{x}}|^\tau = d^\tau$  and

$$B(\mathbf{x}^{\otimes \tau}, \mathbf{y}^{\otimes \tau}) = |\tilde{\mathbf{x}} \cap \tilde{\mathbf{y}}|^\tau / |\tilde{\mathbf{x}}|^\tau = (1 - r/d)^\tau$$

where  $r = \|\mathbf{x} - \mathbf{y}\|_1$ . For every choice of constants  $0 < b_2 < b_1 < 1$  we can choose  $d, \tau, r$ , and  $c \geq 1$  such that  $(1 - r/d)^\tau \approx b_1$  and  $(1 - cr/d)^\tau \approx b_2$ . Now, if there existed an LSH family for Braun-Blanquet with  $\rho < \log(1/b_1)/\log(1/b_2)$  we would be able to obtain an LSH family for Hamming space with

$$\rho < \log(1/b_1)/\log(1/b_2) = \log(1/(1-r/d))/\log(1/(1-cr/d)) \leq 1/c.$$

For appropriate choices of parameters this would contradict the O'Donnell et al. LSH lower bound of  $\rho \gtrsim 1/c$  for Hamming space. The proof itself is mostly an exercise in setting parameters and applying the right bounds and approximations to make everything fit together with the intuition above. Importantly, we use sampling in order to map to a dimension that is much lower than the  $d^\tau$  from the proof sketch in order to make the proof hold for small values of  $p_2$  in relation to  $d$ .

**Hamming to Braun-Blanquet Similarity.** Let  $d \in \mathbb{N}$  and let  $0 < b_2 < b_1 < 1$  be constant as in Theorem 1.2. Let  $\varepsilon \geq 1/d$  be a parameter to be determined. We want to show how to use a transformation  $T: \{0, 1\}^D \rightarrow \{0, 1\}^d$  from Hamming distance to Braun-Blanquet similarity together with our family  $\mathcal{H}_B$  to construct a  $(r, cr, p'_1, p'_2)$ -sensitive family  $\mathcal{H}_H$  for  $D$ -dimensional Hamming space with parameters

$$\begin{aligned} D &= 2^d \\ r &= \sqrt{D} \\ c &= \frac{\ln(1/(b_2 - \varepsilon))}{\ln(1/(b_1 + \varepsilon))} \end{aligned}$$

where  $p'_1$  and  $p'_2$  remain to be determined.

The function  $T$  takes as parameters positive integers  $t, l$ , and  $\tau$ . The output of  $T$  consists of  $t$  concatenated  $l$ -bit strings, each of Hamming weight one. Each of the  $t$  strings is constructed independently at random according to the following process: Sample a vector of indices  $\mathbf{i} = (i_1, i_2, \dots, i_\tau)$  uniformly at random from  $[D]^\tau$  and define  $\mathbf{x}_i \in \{0, 1\}^\tau$  as  $\mathbf{x}_i = \mathbf{x}_{i_1} \circ \mathbf{x}_{i_2} \circ \dots \circ \mathbf{x}_{i_\tau}$ . Let  $\mathbf{z}(\mathbf{x}) \in \{0, 1\}^{2^\tau}$  be indexed by  $j \in \{0, 1\}^\tau$  and set the bits of  $\mathbf{z}(\mathbf{x})$  as follows:

$$\mathbf{z}(\mathbf{x})_j = \begin{cases} 1 & \text{if } \mathbf{x}_i = j, \\ 0 & \text{otherwise.} \end{cases}$$

Next we apply a random function  $g: \{0, 1\}^\tau \rightarrow [l]$  in order to map  $\mathbf{z}(\mathbf{x})$  down to an  $l$ -bit string  $\mathbf{r}(\mathbf{z}(\mathbf{x}))$  of Hamming weight one while

approximately preserving Braun-Blanquet similarity. For  $i \in [l]$  we set

$$\mathbf{r}(\mathbf{z}(\mathbf{x}))_i = \bigvee_{j: g(j)=i} \mathbf{z}(\mathbf{x})_j.$$

Finally we set

$$T(\mathbf{x}) = \mathbf{r}_1(\mathbf{z}_1(\mathbf{x})) \circ \mathbf{r}_2(\mathbf{z}_2(\mathbf{x})) \circ \dots \circ \mathbf{r}_t(\mathbf{z}_t(\mathbf{x}))$$

where each  $\mathbf{r}_i(\mathbf{z}_i(\mathbf{x}))$  is constructed independently at random.

We state the properties of  $T$  for the following parameter setting:

$$\begin{aligned} \tau &= \lfloor \sqrt{D} \ln(1/(b_1 + \varepsilon)) \rfloor \\ l &= \lceil 8/\varepsilon \rceil \\ t &= \lfloor d/l \rfloor. \end{aligned}$$

**LEMMA 4.4.** *For every  $d \in \mathbb{N}$  and  $D = 2^d$  there exists a distribution over functions of the form  $T: \{0, 1\}^D \rightarrow \{0, 1\}^d$  such that for all  $\mathbf{x}, \mathbf{y} \in \{0, 1\}^D$  and random  $T$ :*

- (1)  $|T(\mathbf{x})| = t$ .
- (2) If  $\|\mathbf{x} - \mathbf{y}\|_1 \leq r$  then  $B(T(\mathbf{x}), T(\mathbf{y})) \geq b_1$  with probability at least  $1 - e^{-t\varepsilon^2/2}$ .
- (3) If  $\|\mathbf{x} - \mathbf{y}\|_1 > cr$  then  $B(T(\mathbf{x}), T(\mathbf{y})) < b_2$  with probability at least  $1 - 2e^{t\varepsilon^2/32}$ .

**PROOF.** The first property is trivial. For the second property we consider  $\mathbf{x}, \mathbf{y}$  with  $\|\mathbf{x} - \mathbf{y}\|_1 \leq r$  where we would like to lower bound

$$B(T(\mathbf{x}), T(\mathbf{y})) = \frac{|T(\mathbf{x}) \cap T(\mathbf{y})|}{\max(|T(\mathbf{x})|, |T(\mathbf{y})|)}.$$

We know that  $|T(\mathbf{x})| = |T(\mathbf{y})| = t$  so it remains to lower bound the size of the intersection  $|T(\mathbf{x}) \cap T(\mathbf{y})|$ . Consider the expectation

$$\mathbb{E}[|T(\mathbf{x}) \cap T(\mathbf{y})|] = t \Pr[\mathbf{z}(\mathbf{x}) = \mathbf{z}(\mathbf{y})].$$

We have that  $\mathbf{z}(\mathbf{x}) = \mathbf{z}(\mathbf{y})$  if  $\mathbf{x}$  and  $\mathbf{y}$  take on the same value in the  $\tau$  underlying bit-positions that are sampled to construct  $\mathbf{z}$ . Under the assumption that  $\varepsilon \geq 1/d$ , then for  $d$  greater than some sufficiently large constant we can use a standard approximation to the exponential function (detailed in Lemma A.4 in Appendix A) to show that

$$\begin{aligned} \Pr[\mathbf{z}(\mathbf{x}) = \mathbf{z}(\mathbf{y})] &\geq (1 - r/D)^\tau \\ &\geq (1 - 1/\sqrt{D})^{\sqrt{D} \ln(1/(b_1 + \varepsilon))} \\ &\geq e^{\ln(b_1 + \varepsilon)} (1 - (\ln(b_1 + \varepsilon))^2 / \sqrt{D}) \\ &\geq b_1 + \varepsilon/2. \end{aligned}$$

Seeing as  $|T(\mathbf{x}) \cap T(\mathbf{y})|$  is the sum of  $t$  independent Bernoulli trials we can apply Hoeffding's inequality to yield the following bound:

$$\Pr[|T(\mathbf{x}) \cap T(\mathbf{y})| \leq b_1 t] \leq e^{-t\varepsilon^2/2}.$$

This proves the second property of  $T$ .

For the third property we consider the Braun-Blanquet similarity of distant pairs of points  $\mathbf{x}, \mathbf{y}$  with  $\|\mathbf{x} - \mathbf{y}\|_1 > cr$ . Again, under our assumption that  $\varepsilon \geq 1/d$  and for  $d$  greater than some constant we

have

$$\begin{aligned}
\Pr[\mathbf{z}(\mathbf{x}) = \mathbf{z}(\mathbf{y})] &\leq (1 - cr/D)^\tau \\
&\leq \frac{\left(1 - \frac{\ln(1/(b_2 - \varepsilon))}{\sqrt{D} \ln(1/(b_1 + \varepsilon))}\right)^{\sqrt{D} \ln(1/(b_1 + \varepsilon))}}{1 - c/\sqrt{D}} \\
&\leq (1 + 2c/\sqrt{D})(b_2 - \varepsilon) \\
&\leq b_2 - \varepsilon/2.
\end{aligned}$$

There are two things that can cause the event  $B(T(\mathbf{x}), T(\mathbf{y})) < b_2$  to fail. First, the sum of the  $t$  independent Bernoulli trials for the event  $\mathbf{z}(\mathbf{x}) = \mathbf{z}(\mathbf{x}')$  can deviate too much from its expected value. Second, the mapping down to  $l$ -bit strings that takes place from  $\mathbf{z}(\mathbf{x})$  to  $\mathbf{r}(\mathbf{z}(\mathbf{x}))$  can lead to an additional increase in the similarity due to collisions. Let  $Z$  denote the sum of the  $t$  Bernoulli trials for the events  $\mathbf{z}(\mathbf{x}) = \mathbf{z}(\mathbf{x}')$  associated with  $T$ . We again apply a standard Hoeffding bound to show that

$$\Pr[Z \geq (b_2 - \varepsilon/4)t] \leq e^{-t\varepsilon^2/8}.$$

Let  $X$  denote the number of collisions when performing the universe reduction to  $l$ -bit strings. By our choice of  $l$  we have that  $E[X] \leq (\varepsilon/8)t$ . Another application of Hoeffding's inequality shows that

$$\Pr[X \geq (\varepsilon/4)t] \leq e^{-t\varepsilon^2/32}.$$

We therefore get that

$$\Pr[|T(\mathbf{x}) \cap T(\mathbf{x}')| \geq b_2 t] \leq 2e^{-t\varepsilon^2/32}.$$

This proves the third property of  $T$ .  $\square$

**Contradiction.** To summarize, using the random map  $T$  together with the LSH family  $\mathcal{H}_B$  we can obtain an  $(r, cr, p_1', p_2')$ -sensitive family  $\mathcal{H}_H$  for  $D$ -dimensional Hamming space with  $p_1' = p_1 - \delta$  and  $p_2' = p_2 + \delta$  for  $\delta = 2e^{-t\varepsilon^2/32}$ . For our choice of  $c = \frac{\ln(1/(b_2 - \varepsilon))}{\ln(1/(b_1 + \varepsilon))}$  we plug the family  $\mathcal{H}_H$  into the lower bound of Lemma 4.3 and use that  $O(D^{-1/4}) = O(\varepsilon)$  which follows from our constraint that  $\varepsilon \geq 1/d$ .

$$\begin{aligned}
\rho(\mathcal{H}_H) &\geq 1/c - O(D^{-1/4}) \\
&= \frac{\ln(1/(1 + \varepsilon/b_1)) + \ln(1/b_1)}{\ln(1/(1 - \varepsilon/b_2)) + \ln(1/b_2)} - O(\varepsilon) \\
&\geq \frac{\ln(1/b_1) - \varepsilon/b_1}{\ln(1/b_2) + 2\varepsilon/b_2} - O(\varepsilon) \\
&= \frac{\ln(1/b_1)}{\ln(1/b_2)} - O(\varepsilon)
\end{aligned}$$

Under our assumed properties of  $\mathcal{H}_B$ , we can upper bound the value of  $\rho$  for  $\mathcal{H}_H$ . For simplicity we temporarily define  $\lambda = 2\delta/p_2$  and assume that  $\lambda/\ln(1/p_2) \leq 1/2$  and  $\ln(1/p_2) \geq 1$ . The latter property holds without loss of generality through use of the standard LSH powering technique [21, 23, 30] that allows us to transform an LSH family with  $p_2 < 1$  to a family that has  $p_2 \leq 1/e$  without changing

its associated  $\rho$ -value.

$$\begin{aligned}
\rho(\mathcal{H}_H) &= \frac{\ln(1/p_1')}{\ln(1/p_2')} = \frac{\ln(1/p_1) + \ln(1/(1 - \delta/p_1))}{\ln(1/p_2) + \ln(1/(1 + \delta/p_2))} \\
&\leq \frac{\ln(1/p_1) + \lambda}{\ln(1/p_2) - \lambda} = \frac{\ln(1/p_1) + \lambda}{(\ln(1/p_2))(1 - \lambda/(\ln(1/p_2)))} \\
&\leq \frac{\ln(1/p_1) + \lambda}{\ln(1/p_2)} (1 + 2\lambda/(\ln(1/p_2))) = \frac{\ln(1/p_1)}{\ln(1/p_2)} + O(\delta/p_2) \\
&\leq \frac{\ln(1/b_1)}{\ln(1/b_2)} - \gamma + O(\delta/p_2).
\end{aligned}$$

We get a contradiction between our upper bound and lower bound for  $\rho(\mathcal{H}_H)$  whenever  $\gamma$  violates the following relation that summarizes the bounds:

$$\frac{\ln(1/b_1)}{\ln(1/b_2)} - O(\varepsilon) \leq \rho(\mathcal{H}_H) \leq \frac{\ln(1/b_1)}{\ln(1/b_2)} - \gamma + O(\delta/p_2).$$

In order for a contradiction to occur, the value of  $\gamma$  has to satisfy

$$\gamma > O(\varepsilon) + O(\delta/p_2).$$

By our setting of  $t = \lfloor d/l \rfloor$  and  $l = \lceil 8/\varepsilon \rceil$  we have that  $\delta = e^{-\Omega(d\varepsilon^3)}$ . We can cause a contradiction for a setting of  $\varepsilon^3 = K \frac{\ln(d/p_2)}{d}$  where  $K$  is some constant and where we assume that  $d$  is greater than some constant. The value of  $\gamma$  for which the lower bound holds can be upper bounded by

$$\gamma = O\left(\frac{\ln(d/p_2)}{d}\right)^{1/3}.$$

This completes the proof of Theorem 1.2.

## 5 EQUIVALENT SET SIMILARITY PROBLEMS

In this section we consider how to use our data structure for Braun-Blanquet similarity search to support other similarity measures such as Jaccard similarity. We already observed in the introduction that a direct translation exists between several similarity measures whenever the size of every sets is fixed to  $t$ . Call an  $(s_1, s_2)$ - $S$ -similarity search problem  $(t, t')$ -regular if  $P$  is restricted to vectors of weight  $t$  and queries are restricted to vectors of weight  $t'$ . Obviously, a  $(t, t')$ -regular similarity search problem is no harder than the general similarity search problem, but it also cannot be too much easier when expressed as a function of the thresholds  $(s_1, s_2)$ : For every pair  $(t, t') \in \{0, \dots, d\}^2$  we can construct a  $(t, t')$ -regular data structure (such that each point  $\mathbf{x} \in P$  is represented in the  $d + 1$  data structures with  $t = |\mathbf{x}|$ ), and answer a query for  $\mathbf{q} \in \{0, 1\}^d$  by querying all data structures with  $t' = |\mathbf{q}|$ . Thus, the time and space for the general  $(s_1, s_2)$ - $S$ -similarity search problem is at most  $d + 1$  times larger than the time and space of the most expensive  $(t, t')$ -regular data structure. This does *not* mean that we cannot get better bounds in terms of other parameters, and in particular we expect that  $(t, t')$ -regular similarity search problems have difficulty that depends on parameters  $t$  and  $t'$ .

**Dimension Reduction.** If the dimension is large a factor of  $d$  may be significant. However, for most natural similarity measures a  $(s_1, s_2)$ - $S$ -similarity problem in  $d \gg (\log n)^3$  dimensions can be reduced to a logarithmic number of  $(s'_1, s'_2)$ - $S$ -similarity problems on  $P' \subseteq \{0, 1\}^{d'}$  in  $d' = (\log n)^3$  dimensions with  $s'_1 = s_1 - O(1/\log n)$  and  $s'_2 = s_2 + O(1/\log n)$ . Since the similarity gap is close to the one

in the original problem,  $s'_1 - s'_2 = s_1 - s_2 - O(1/\log n)$ , where  $s_1$  and  $s_2$  are assumed to be independent of  $n$ , the difficulty ( $\rho$ -value) remains essentially the same. First, split  $P$  into  $\log d$  size classes  $P_i$  such that vectors in class  $i$  have size in  $[2^i; 2^{i+1})$ . For each size class the reduction is done independently and works by a standard technique: sample a sequence of random sets  $I_j \subseteq \{1, \dots, d\}$ ,  $i = 1, \dots, d'$ , and set  $\mathbf{x}'_j = \vee_{\ell \in I_j} \mathbf{x}_\ell$ . The size of each set  $I_j$  is chosen such that  $\Pr[\mathbf{x}'_j = 1] \approx 1/\log(n)$  when  $|\mathbf{x}| = 2^{i+1}$ . By Chernoff bounds this mapping preserves the relative weight of vectors up to size  $2^i \log n$  up to an additive  $O(1/\log n)$  term with high probability. Assume now that the similarity measure  $S$  is such that for vectors in  $P_i$  we only need to consider  $|\mathbf{q}|$  in the range from  $2^i/\log n$  to  $2^i \log n$  (since if the size difference is larger, the similarity is negligible). The we can apply Chernoff bounds to the relative weights of the dimension-reduced vectors  $\mathbf{x}'$ ,  $\mathbf{q}'$  and the intersection  $\mathbf{x}' \cap \mathbf{q}'$ . In particular, we get that the Jaccard similarity of a pair of vectors is preserved up to an additive error of  $O(1/\log n)$  with high probability. The class of similarity measures for which dimension reduction to  $(\log n)^{O(1)}$  dimensions is possible is large, and we do not attempt to characterize it here. Instead, we just note that for such similarity measures we can determine the complexity of similarity search up to a factor  $(\log n)^{O(1)}$  by only considering regular search problems.

**Equivalence of Regular Similarity Search Problems.** We call a set similarity measure on  $\{0, 1\}^d$  *symmetric* if it can be written in the form  $S(\mathbf{q}, \mathbf{x}) = f_{d, |\mathbf{q}|, |\mathbf{x}|}(|\mathbf{q} \cap \mathbf{x}|)$ , where each function  $f_{d, |\mathbf{q}|, |\mathbf{x}|}: \mathbb{N} \rightarrow [0; 1]$  is nondecreasing. All 59 set similarity measures listed in the survey [15], normalized to yield similarities in  $[0; 1]$ , are symmetric. In particular this is the case for Jaccard similarity (where  $J(\mathbf{q}, \mathbf{x}) = |\mathbf{q} \cap \mathbf{x}| / (|\mathbf{q}| + |\mathbf{x}| - |\mathbf{q} \cap \mathbf{x}|)$ ) and for Braun-Blanquet similarity. For a symmetric similarity measure  $S$ , the predicate  $S(\mathbf{q}, \mathbf{x}) \geq s_1$  is equivalent to the predicate  $|\mathbf{q} \cap \mathbf{x}| \geq i_1$ , where  $i_1 = \min\{i \mid f_{d, t', t}(i) \geq s_1\}$ , and  $S(\mathbf{q}, \mathbf{x}) > s_2$  is equivalent to the predicate  $|\mathbf{q} \cap \mathbf{x}| \geq i_2$ , where  $i_2 = \min\{i \mid f_{d, t', t}(i) > s_2\}$ . This means that every  $(t, t')$ -regular  $(s_1, s_2)$ - $S$ -similarity search problem on  $P \subseteq \{0, 1\}^d$  is equivalent to an  $(i_1/d, i_2/d)$ - $I$ -similarity search problem on  $P$ , where  $I(\mathbf{q}, \mathbf{x}) = |\mathbf{x} \cap \mathbf{q}|/d$ . In other words, all symmetric similarity search problems can be translated to each other, and it suffices to study a single one, such as Braun-Blanquet similarity.

**Jaccard similarity.** We briefly discuss Jaccard similarity since it is the most widely used measure of set similarity. If we consider the problem of  $(j_1, j_2)$ -approximate Jaccard similarity search in the  $(t, t')$ -regular case with  $t \neq t'$  then our Theorem 1.1 is no longer guaranteed to yield the lowest value of  $\rho$  among competing data-independent approaches such as MinHash and Angular LSH. To simplify the comparison between different measures we introduce parameters  $\beta$  and  $b$  defined by  $|y| = \beta|x|$  and  $b = |\mathbf{x} \cap \mathbf{y}|/|\mathbf{x}|$  (note that  $0 \leq b \leq \beta \leq 1$ ). The three primary measures of set similarity considered in this paper can then be written as follows:

$$\begin{aligned} B(\mathbf{x}, \mathbf{y}) &= b \\ J(\mathbf{x}, \mathbf{y}) &= \frac{b}{1 + \beta - b} \\ C(\mathbf{x}, \mathbf{y}) &= \frac{b}{\sqrt{\beta}} \end{aligned}$$

As shown in Figure 6 among angular LSH, MinHash, and CHOSEN PATH, the technique with the lowest  $\rho$ -value is different depending on the parameters  $(j_1, j_2)$  and asymmetry  $\beta$ . We know that CHOSEN PATH is optimal and strictly better than the competing data-independent techniques across the entire parameter space  $(j_1, j_2)$  when  $\beta = 1$ , but it remains open to find tight upper and lower bounds in the case where  $\beta \neq 1$ .

## 6 CONCLUSION AND OPEN PROBLEMS

We have seen that, perhaps surprisingly, there exists a relatively simple way of strictly improving the  $\rho$ -value for data-independent set similarity search in the case where all sets have the same size. To implement the required locality-sensitive map efficiently we introduce a new technique based on branching processes that could possibly lead to more efficient solutions in other settings.

It remains an open problem to find tight upper and lower bounds on the  $\rho$ -value for Jaccard and cosine similarity search that hold for the entire parameter space in the general setting with arbitrary set sizes. Perhaps a modified version of the CHOSEN PATH algorithm can yield an improved solution to Jaccard similarity search in general. One approach is to generalize the condition  $h_i(p \circ j) < \mathbf{x}_j/b_1|\mathbf{x}|$  to use different thresholds for queries and updates. This yields different space-time tradeoffs when applying the CHOSEN PATH algorithm to Jaccard similarity search.

Another interesting question is if the improvement shown for sparse vectors can be achieved in general for inner product similarity. A similar, but possibly easier, direction would be to consider *weighted* Jaccard similarity.

## ACKNOWLEDGMENTS

We thank Thomas Dybdahl Ahle for comments on a previous version of this manuscript.

## A DETAILS BEHIND THE LOWER BOUND

### A.1 Tools

For clarity we state some standard technical lemmas that we use to derive LSH lower bounds.

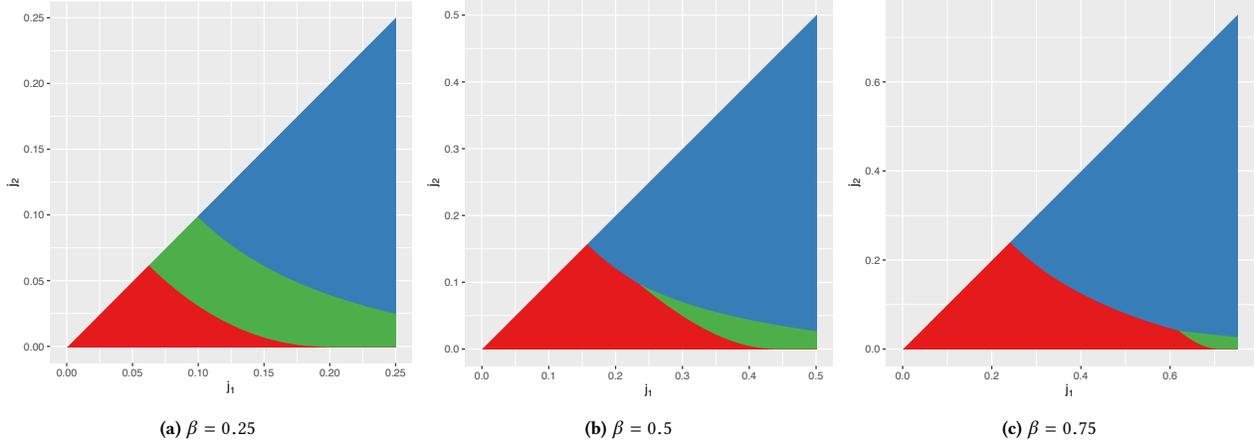
**LEMMA A.1 (HOEFFDING [22, THEOREM 1]).** *Let  $X_1, X_2, \dots, X_n$  be independent random variables satisfying  $0 \leq X_i \leq 1$  for  $i \in [n]$ . Define  $X = X_1 + X_2 + \dots + X_n$ ,  $Z = X/n$ , and  $\mu = E[Z]$ , then:*

- For  $\hat{\mu} \geq \mu$  and  $0 < \varepsilon < 1 - \hat{\mu}$  we have that  $\Pr[Z - \hat{\mu} \geq \varepsilon] \leq e^{-2n\varepsilon^2}$ .
- For  $\hat{\mu} \leq \mu$  and  $0 < \varepsilon < \hat{\mu}$  we have that  $\Pr[Z - \hat{\mu} \leq -\varepsilon] \leq e^{-2n\varepsilon^2}$ .

**LEMMA A.2 (CHERNOFF [27, THM. 4.4 AND 4.5]).** *Let  $X_1, \dots, X_n$  be independent Poisson trials and define  $X = \sum_{i=1}^n X_i$  and  $\mu = E[X]$ . Then, for  $0 < \varepsilon < 1$  we have*

- $\Pr[X \geq (1 + \varepsilon)\mu] \leq e^{-\varepsilon^2\mu/3}$ .
- $\Pr[X \leq (1 - \varepsilon)\mu] \leq e^{-\varepsilon^2\mu/2}$ .

**LEMMA A.3 (BOUNDING THE LOGARITHM [36]).** *For  $x > -1$  we have that  $\frac{x}{1+x} \leq \ln(1+x) \leq x$ .*



**Figure 6: Solution with lowest  $\rho$ -value for the  $(j_1, j_2)$ -approximate Jaccard similarity search problem for different values of  $\beta$ . Blue is angular LSH. Green is MinHash. Red is CHOSEN PATH. Note the difference in the axes for different values of  $\beta$  as it must hold that  $0 \leq j_2 \leq j_1 \leq \beta$ .**

LEMMA A.4 (APPROXIMATING THE EXPONENTIAL FUNCTION [29, PROP. B.3]). *For all  $t, n \in \mathbb{R}$  with  $|t| \leq n$  we have that  $e^t(1 - \frac{t^2}{n}) \leq (1 + \frac{t}{n})^n \leq e^t$ .*

## A.2 Proof of Lemma 4.3

**Preliminaries.** We will reuse the notation of Section 3. from O’Donnell et al. [30].

*Definition A.5.* For  $0 \leq \lambda < 1$  we say that  $(\mathbf{x}, \mathbf{y})$  are  $(1 - \lambda)$ -correlated if  $\mathbf{x}$  is chosen uniformly at random from  $\{0, 1\}^d$  and  $\mathbf{y}$  is constructed by rerandomizing each bit from  $\mathbf{x}$  independently at random with probability  $\lambda$ .

Let  $(\mathbf{x}, \mathbf{y})$  be  $e^{-t}$ -correlated and let  $\mathcal{H}$  be a family of hash functions on  $\{0, 1\}^d$ , then we define

$$\mathbb{K}_{\mathcal{H}}(t) = \Pr_{\substack{h \sim \mathcal{H} \\ (\mathbf{x}, \mathbf{y}) e^{-t}\text{-corr}^d}} [h(\mathbf{x}) = h(\mathbf{y})].$$

We have that  $\mathbb{K}_{\mathcal{H}}(t)$  is a log-convex function which implies the following property that underlies the lower bound:

LEMMA A.6. *For every family of hash functions  $\mathcal{H}$  on  $\{0, 1\}^d$ , every  $t \geq 0$ , and  $c \geq 1$  we have*

$$\frac{\ln(1/\mathbb{K}_{\mathcal{H}}(t))}{\ln(1/\mathbb{K}_{\mathcal{H}}(ct))} \geq \frac{1}{c}. \quad (6)$$

The idea behind the proof is to tie  $p_1$  to  $\mathbb{K}_{\mathcal{H}}(t)$  and  $p_2$  to  $\mathbb{K}_{\mathcal{H}}(ct)$  through Chernoff bounds and then apply Lemma A.6 to show that  $\rho \gtrsim 1/c$ .

**Proof.** Begin by assuming that we have a family  $\mathcal{H}$  that satisfies the conditions of Lemma 4.3. Note that the expected Hamming distance between  $(1 - \lambda)$ -correlated points  $\mathbf{x}$  and  $\mathbf{y}$  is given by  $(\lambda/2)d$ . We set  $\lambda_{p_1}/2 = d^{-1/2} - d^{-5/8}$  and  $\lambda_{p_2}/2 = cd^{-1/2} + 2cd^{-5/8}$  and let  $(\mathbf{x}, \mathbf{y})$  denote  $(1 - \lambda_{p_1})$ -correlated random strings and  $(\mathbf{x}, \mathbf{x}')$  denote  $(1 - \lambda_{p_2}q)$ -correlated random strings. By standard Chernoff bounds

we get the following guarantees:

$$\begin{aligned} \Pr[\|\mathbf{x} - \mathbf{y}\|_1 \geq r] &\leq e^{-\Omega(d^{1/4})}, \\ \Pr[\|\mathbf{x} - \mathbf{x}'\|_1 \leq cr] &\leq e^{-\Omega(d^{1/4})}. \end{aligned}$$

We will establish a relationship between  $\mathbb{K}_{\mathcal{H}}(t_{p_1})$  and  $p_1$  on the one hand, and  $\mathbb{K}_{\mathcal{H}}(t_{p_2})$  and  $p_2$  on the other hand, for the following choice of parameters  $t_{p_1}$  and  $t_{p_2}$ :

$$\begin{aligned} t_{p_1} &= -\ln(1 - 2(d^{-1/2} - d^{-5/8})) \\ t_{p_2} &= -\ln(1 - 2c(d^{-1/2} + 2d^{-5/8})). \end{aligned}$$

By the properties of  $\mathcal{H}$  and from the definition of  $\mathbb{K}_{\mathcal{H}}$  we have that

$$\begin{aligned} \mathbb{K}_{\mathcal{H}}(t_{p_1}) &\geq p_1(1 - \Pr[\|\mathbf{x} - \mathbf{y}\|_1 > r]) \geq p_1 - \Pr[\|\mathbf{x} - \mathbf{y}\|_1 \geq r] \\ \mathbb{K}_{\mathcal{H}}(t_{p_2}) &\leq p_2(1 - \Pr[\|\mathbf{x} - \mathbf{x}'\|_1 \leq cr]) + \Pr[\|\mathbf{x} - \mathbf{x}'\|_1 \leq cr] \\ &\leq p_2 + \Pr[\|\mathbf{x} - \mathbf{x}'\|_1 \leq cr]. \end{aligned}$$

Let  $\delta = \max\{\Pr[\|\mathbf{x} - \mathbf{y}\|_1 \geq r], \Pr[\|\mathbf{x} - \mathbf{x}'\|_1 \leq cr]\} = e^{-\Omega(d^{1/4})}$ . By Lemma A.6 and our setting of  $t_{p_1}$  and  $t_{p_2}$  we can use the bounds on the natural logarithm from Lemma A.3 to show the following:

$$\begin{aligned} \frac{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_1}))}{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2}))} &\geq \frac{t_{p_1}}{t_{p_2}} = \frac{\ln(1 - 2(d^{-1/2} - d^{-5/8}))}{\ln(1 - 2c(d^{-1/2} + 2d^{-5/8}))} \\ &\geq \frac{2(d^{-1/2} - d^{-5/8})}{2c(d^{-1/2} + 2d^{-5/8})} - 2(d^{-1/2} - d^{-5/8}) \\ &\geq \frac{1 - d^{-1/4}}{c + 2d^{-1/4}} - 2(d^{-1/2} - d^{-5/8}) \\ &= \frac{1}{c} - O(d^{-1/4}). \end{aligned}$$

We proceed by lower bounding  $\rho$  where we make use of the inequalities derived above.

$$\mathbb{K}_{\mathcal{H}}(t_{p_2}) - \delta \leq p_2 < p_1 \leq \mathbb{K}_{\mathcal{H}}(t_{p_1}) + \delta.$$

By Lemma A.6 combined with the restrictions on our parameters, for  $d$  greater than some constant we have that  $\mathbb{K}_{\mathcal{H}}(t_{p_2}) \geq$

$\mathbb{K}_{\mathcal{H}}(t_{p_1})^{2c} \geq (p_1/2)^{2c} \geq (2d)^{-2c} \geq (2d)^{-2d^{1/8}}$ . Furthermore, we lower bound  $\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2}))$  by using that  $\mathbb{K}_{\mathcal{H}}(t_{p_2}) \leq p_2 + \delta$  together with the restriction that  $p_2 \geq 1 - 1/d$  and the properties of  $\delta$ . For  $d$  greater than some constant it therefore holds that  $\mathbb{K}_{\mathcal{H}}(t_{p_2}) \leq 1 - 1/2d$  from which it follows that  $\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2})) \geq 1/2d$ .

$$\begin{aligned} \frac{\ln(1/p_1)}{\ln(1/p_2)} &\geq \frac{\ln(1/(\mathbb{K}_{\mathcal{H}}(t_{p_1}) + \delta))}{\ln(1/(\mathbb{K}_{\mathcal{H}}(t_{p_2}) - \delta))} \\ &= \frac{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_1})) - \ln(1 + \delta/\mathbb{K}_{\mathcal{H}}(t_{p_1}))}{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2})) + \ln(1/(1 - \delta/\mathbb{K}_{\mathcal{H}}(t_{p_2})))} \\ &\geq \frac{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_1})) - \delta/\mathbb{K}_{\mathcal{H}}(t_{p_1})}{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2})) + 2\delta/\mathbb{K}_{\mathcal{H}}(t_{p_2})} \\ &\geq \frac{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_1}))}{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2}))} - \frac{3\delta}{\mathbb{K}_{\mathcal{H}}(t_{p_2}) \ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2}))}. \end{aligned}$$

By the arguments above we have that

$$\frac{3\delta}{\mathbb{K}_{\mathcal{H}}(t_{p_2}) \ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2}))} = e^{-\Omega(d^{1/4})} = O(d^{-1/4}).$$

Inserting the lower bound for  $\frac{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_1}))}{\ln(1/\mathbb{K}_{\mathcal{H}}(t_{p_2}))}$  results in the lemma.

## B COMPARISONS

For completeness we state the proofs behind the comparisons between the  $\rho$ -values obtained by the CHOSEN PATH algorithm and other LSH techniques.

### B.1 MinHash

For data sets with fixed sparsity and Braun-Blanquet similarities  $0 < b_2 < b_1 < 1$  we have that  $\rho/\rho_{\text{minhash}} = f(b_2)/f(b_1)$  where  $f(x) = \log(x/(2-x))/\log(x)$ . If  $f(x)$  is monotone increasing in  $(0; 1)$  then  $\rho/\rho_{\text{minhash}} < 1$ . For  $x \in (0; 1)$  we have that  $\text{sign}(f'(x)) = \text{sign}(g(x))$  where  $g(x) = \ln(x) + (2-x)\ln(2-x)$ . The function  $g(x)$  equals zero at  $x = 1$  and has the derivative  $g'(x) = \ln(x) - \ln(2-x)$  which is negative for values of  $x \in (0; 1)$ . We can therefore see that  $f'(x)$  is positive in the interval and it follows that  $\rho < \rho_{\text{minhash}}$  for every choice of  $0 < b_2 < b_1 < 1$ .

### B.2 Angular LSH

We have that  $\rho/\rho_{\text{angular}} < 1$  if  $f(x) = \ln(x)\frac{1+x}{1-x}$  is a monotone increasing function for  $x \in (0; 1)$ . For  $x \in (0; 1)$  we have that  $\text{sign}(f'(x)) = \text{sign}(g(x))$  where  $g(x) = (1-x^2)/2 + x \ln x$ . We note that  $g(1) = 0$  and  $g'(x) = 1 - x + \ln x$ . Therefore, if  $g'(x) < 0$  for  $x \in (0; 1)$  it holds that  $g(x) > 0$  and  $f(x)$  is monotone increasing in the same interval. We have that  $g'(1) = 0$  and  $g''(x) = -1 + 1/x > 0$  implying that  $g'(x) < 0$  in the interval.

### B.3 Data-dependent LSH

LEMMA B.1. *Let  $0 < b_2 < b_1 < 1$  and fix  $\rho = 1/2$  such that  $b_1 = \sqrt{b_2}$ . Then we have that  $\rho < \rho_{\text{datadep}}$  for every value of  $b_2 < 1/4$ .*

PROOF. We will compare  $\rho = \log(b_1)/\log(b_2)$  and  $\rho_{\text{datadep}} = \frac{1-b_1}{1+b_1-2b_2}$  when  $\rho$  is fixed at  $\rho = 1/2$ , or equivalently,  $b_1 = \sqrt{b_2}$ . We can solve the quadratic equation  $1/2 = \frac{1-\sqrt{b_2}}{1+\sqrt{b_2}-2b_2}$  to see that for  $0 < b_2 < 1$  we have that  $\rho = \rho_{\text{datadep}}$  only when  $b_2 = 1/4$ . The derivative of  $\rho_{\text{datadep}}$  with respect to  $b_2$  is negative when  $b_1 = \sqrt{b_2}$ .

Under this restriction we therefore have that  $\rho < \rho_{\text{datadep}}$  for  $b_2 < 1/4$  which is equivalent to  $j_2 < 1/7$  in the fixed-weight setting.  $\square$

To compare  $\rho$ -values over the full parameter space we use the following two lemmas.

LEMMA B.2. *For every choice of fixed  $0 < \rho < 1$  let  $b_2 = b_1^{1/\rho}$ . Then  $\rho_{\text{datadep}} = \frac{1-b_1}{1+b_1-2b_2}$  is decreasing in  $b_1$  for  $b_1 \in (0; 1)$ .*

PROOF. The sign of the derivative of  $\rho_{\text{datadep}}$  with respect to  $b_1$  is equal to the sign of the function  $g(x) = -\rho x^{-1/\rho} + \rho - 1 + x^{-1}$  for  $x \in (0; 1)$ . We have that  $g(1) = 0$  and  $g'(x) = x^{-1/\rho} - 1 - x^{-2} > 0$  for  $x \in (0; 1)$  which shows that  $g(x) < 0$  in the interval.  $\square$

LEMMA B.3. *For  $1/5 = b_2 < b_1 < 1$  we have that  $\rho < \rho_{\text{datadep}}$ .*

PROOF. For fixed  $b_2 = 1/5$  consider  $f(b_1) = \rho - \rho_{\text{datadep}}$  as a function of  $b_1$  in the interval  $[1/5, 1]$ . We want to show that  $f(b_1) < 0$  for  $b_1 \in (1/5; 1)$ . In the endpoints the function takes the value 0. Between the endpoints we find that  $f'(b_1) = \frac{1}{\ln(5)b_1} + \frac{8/5}{(3/5+b_1)^2}$  and that  $f'(b_1) = 0$  is a quadratic form with only one solution  $b_1^*$  in  $[1/5; 1]$ . By Lemma B.1 we know that for  $b_2 = 1/5$  and  $b_1 = 1/\sqrt{5}$  it holds that  $f(b_1) < 0$ . Since  $f(1/5) = f(1) = 0$ ,  $f'(b_1) = 0$  only in a single point in  $[1/5; 1]$ , and  $f(1/\sqrt{5}) < 0$  we can conclude that the lemma holds.  $\square$

COROLLARY B.4. *For every choice of  $b_1, b_2$  satisfying  $0 < b_2 \leq 1/5$  and  $b_2 < b_1 < 1$  we have that  $\rho < \rho_{\text{datadep}}$ .*

PROOF. If  $b_2 = 1/5$  the property holds by Lemma B.3. If  $b_2 < 1/5$  we define new variables  $\hat{b}_2, \hat{b}_2$ , setting  $\hat{b}_1 = \hat{b}_1^{\rho(b_1, b_2)}$  and initially consider  $\hat{b}_2 = 1/5$ . In this setting we again have that  $\rho(\hat{b}_1, \hat{b}_2) < \rho_{\text{datadep}}(\hat{b}_1, \hat{b}_2)$ . According to Lemma B.2 it holds that  $\rho_{\text{datadep}}$  is decreasing in  $b_2$  for fixed  $\rho$ . Therefore, as  $\hat{b}_2$  decreases to  $\hat{b}_2 = b_2$  where  $\hat{b}_1 = b_1$  we have that  $\rho(\hat{b}_1, \hat{b}_2) = \rho$  remains constant while  $\rho_{\text{datadep}}$  increases. Since it held that  $\rho < \rho_{\text{datadep}}$  at the initial values of  $\hat{b}_1, \hat{b}_2$  it must also hold for  $b_1, b_2$ .  $\square$

### Numerical Comparison of MinHash and Data-dep. LSH.

Comparing  $\rho_{\text{minhash}}$  to  $\rho_{\text{datadep}}$  we can verify numerically that even for  $b_2$  fixed as low as  $b_2 = 1/23$  we can find values of  $b_1$  (for example  $b_1 = 0.995$  such that  $\rho_{\text{minhash}} > \rho_{\text{datadep}}$ ).

## REFERENCES

- [1] T. D. Ahle, R. Pagh, I. P. Razenshteyn, and F. Silvestri. 2016. On the Complexity of Inner Product Similarity Join. In *Proc. PODS'16*. 151–164.
- [2] A. Andoni, P. Indyk, T. Laarhoven, I. Razenshteyn, and L. Schmidt. 2015. Practical and optimal LSH for angular distance. In *Proc. NIPS '15*. 1225–1233.
- [3] A. Andoni, P. Indyk, H. L. Nguyen, and I. P. Razenshteyn. 2014. Beyond Locality-Sensitive Hashing. In *Proc. SODA '14*. 1018–1028.
- [4] A. Andoni, T. Laarhoven, I. P. Razenshteyn, and E. Waingarten. 2017. Optimal Hashing-based Time-Space Trade-offs for Approximate Near Neighbors. In *Proc. SODA '17*. 47–66.
- [5] A. Andoni and I. Razenshteyn. 2015. Optimal Data-Dependent Hashing for Approximate Near Neighbors. In *Proc. STOC '15*. 793–801.
- [6] A. Andoni and I. Razenshteyn. 2016. Tight Lower Bounds for Data-Dependent Locality-Sensitive Hashing. In *Proc. SoCG '16*. 9:1–9:11.
- [7] Arvind Arasu, Venkatesh Ganti, and Raghav Kaushik. 2006. Efficient exact set-similarity joins. In *Proceedings of the 32nd international conference on Very large data bases*. VLDB Endowment, 918–929.

- [8] Roberto J Bayardo, Yiming Ma, and Ramakrishnan Srikant. 2007. Scaling up all pairs similarity search. In *Proceedings of the 16th international conference on World Wide Web*. ACM, 131–140.
- [9] A. Becker, L. Ducas, N. Gama, and T. Laarhoven. 2016. New directions in nearest neighbor searching with applications to lattice sieving. In *Proc. SODA '16*. 10–24.
- [10] Josias Braun-Blanquet. 1932. *Plant sociology. The study of plant communities*. McGraw-Hill.
- [11] Andrei Z. Broder. 1997. On the resemblance and containment of documents. In *Compression and Complexity of Sequences 1997. Proceedings*. IEEE, 21–29.
- [12] Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. 1997. Syntactic clustering of the web. *Computer Networks and ISDN Systems* 29, 8 (1997), 1157–1166.
- [13] M. Charikar. 2002. Similarity estimation techniques from rounding algorithms. In *Proc. STOC '02*. 380–388.
- [14] F. Chierichetti and R. Kumar. 2015. LSH-Preserving Functions and Their Applications. *J. ACM* 62, 5 (2015), 33.
- [15] S. Choi, S. Cha, and C. C. Tappert. 2010. A survey of binary similarity and distance measures. *J. Syst. Cybern. Informatics* 8, 1 (2010), 43–48.
- [16] T. Christiani. 2017. A Framework for Similarity Search with Space-Time Tradeoffs using Locality-Sensitive Filtering. In *Proc. SODA '17*. 31–46.
- [17] E. Cohen. 1997. Size-estimation framework with applications to transitive closure and reachability. *J. Comp. Syst. Sci.* 55, 3 (1997), 441–453.
- [18] E. Cohen and H. Kaplan. 2009. Leveraging discarded samples for tighter estimation of multiple-set aggregates. *ACM SIGMETRICS Performance Evaluation Review* 37, 1 (2009), 251–262.
- [19] M. Dubiner. 2010. Bucketing coding and information theory for the statistical high-dimensional nearest-neighbor problem. *IEEE Trans. Information Theory* 56, 8 (2010), 4166–4179.
- [20] T. Hagerup. 1998. Sorting and Searching on the Word RAM. In *Proc. STACS '98*. 366–398.
- [21] S. Har-Peled, P. Indyk, and R. Motwani. 2012. Approximate Nearest Neighbor: Towards Removing the Curse of Dimensionality. *Theory of computing* 8, 1 (2012), 321–350.
- [22] W. Hoeffding. 1963. Probability inequalities for sums of bounded random variables. *Jour. Am. Stat. Assoc.* 58, 301 (1963), 13–30.
- [23] P. Indyk and R. Motwani. 1998. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proc. STOC '98*. 604–613.
- [24] T. Laarhoven. 2015. Tradeoffs for nearest neighbors on the sphere. *CoRR* abs/1511.07527 (2015). <http://arxiv.org/abs/1511.07527>
- [25] Ping Li and Arnd Christian König. 2011. Theory and applications of b-bit minwise hashing. *Commun. ACM* 54, 8 (2011), 101–109.
- [26] M. Mitzenmacher, R. Pagh, and N. Pham. 2014. Efficient estimation for high similarities using odd sketches. In *Proc. WWW '14*. 109–118.
- [27] M. Mitzenmacher and E. Upfal. 2005. *Probability and computing*. Cambridge University Press, New York, NY.
- [28] R. Motwani, A. Naor, and R. Panigrahy. 2007. Lower Bounds on Locality Sensitive Hashing. *SIAM J. Discrete Math.* 21, 4 (2007), 930–935.
- [29] Rajeew Motwani and Prabhakar Raghavan. 2010. *Randomized algorithms*. Chapman & Hall/CRC.
- [30] R. O'Donnell, Y. Wu, and Y. Zhou. 2014. Optimal lower bounds for locality-sensitive hashing (except when q is tiny). *ACM Transactions on Computation Theory (TOCT)* 6, 1 (2014), 5.
- [31] Rasmus Pagh, Morten Stöckel, and David P Woodruff. 2014. Is min-wise hashing optimal for summarizing set intersection?. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. ACM, 109–120.
- [32] R. Panigrahy, K. Talwar, and U. Wieder. 2010. Lower Bounds on Near Neighbor Search via Metric Expansion. In *Proc. FOCS '10*. 805–814.
- [33] Anshumali Shrivastava and Ping Li. 2015. Asymmetric minwise hashing for indexing binary inner products and set containment. In *Proceedings of the 24th International Conference on World Wide Web*. ACM, 981–991.
- [34] K. Terasawa and Y. Tanaka. 2007. Spherical LSH for Approximate Nearest Neighbor Search on Unit Hypersphere. In *Proc. WADS '07*. 27–38.
- [35] Mikkel Thorup. 2013. Bottom-k and priority sampling, set similarity and subset sums with minimal independence. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*. ACM, 371–380.
- [36] F. Topsøe. 2007. *Some Bounds for the Logarithmic Function*. Vol. 4. Nova Science, 137–151.
- [37] Albert L Zobrist. 1970. A new hashing method with application for game playing. *ICCA journal* 13, 2 (1970), 69–73.