

Set similarity with b-bit k-permutation Minwise Hashing - DRAFT

Simon Wehrli

28.4.2013

1 Introduction

This report is written for the seminar titled “Algorithms for Database Systems” at ETH Zürich. The seminar participants read and summarize various papers, which treat solving problems in the context of Big Data, which is this year’s topic. This report summarizes and explains the core concepts of [1] and [2].

With Big Data, most problems emphasis shifts from single computational complexity to memory space usage considerations. Typically, the runtime is dominated by the time needed for memory accesses, hence the goal is to reduce the number of memory accesses. Often this is achieved by more compact data structures and at the price of less accurate result.

Many today’s applications are faced with very large datasets. A common task is to find *similarity* between two or several such sets. There are lots of problem solutions which use a mapping to sets of properties and then do a *similarity search* on them. Fast (approximative) algorithms for this search enable improvements of many well-known algorithms, e.g. of machine learning and computer vision.

We start by explaining the original MINWISE HASHING, which nicely presents the basic idea behind all algorithms. We move on to a major improvement, the B-BIT K-PERMUTATION MINWISE HASHING, which reduces storage at the cost of more iterations in the algorithm and is a generalization of the concept. We then present ONE PERMUTATION HASHING, which achieves surprising accuracy using only one permutation. We focus on the concepts and algorithms and will reference to papers for applications in practice.

1.1 Similarity

We denote by Ω the set of all possible items of the sets $S_n \subseteq \Omega$, $n = 1$ to N . $|\Omega| = D$ is always large (e.g. $D = 2^{64}$). Often we consider only two sets $S_1 = X$, $S_2 = Y$. Let $a = |X \cap Y|$, $b = |X| - a$ and $c = |Y| - a$. Figure 1 on the following page summarizes these definitions graphically.

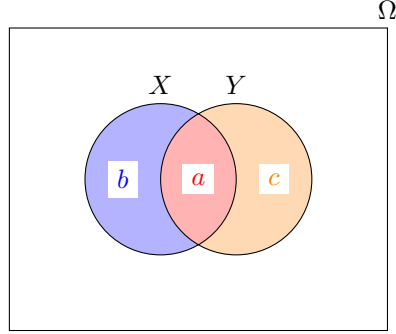


Figure 1: Two example sets in Ω and the notation a, b, c for the sizes of important subsets.

Definition 1. The normalized similarity between two sets X and Y , known as *resemblance* or Jaccard similarity, denoted by R , is

$$R = \frac{|X \cap Y|}{|X \cup Y|} = \frac{a}{|X| + |Y| - a} = \frac{a}{a + b + c} \quad (1)$$

1.2 Motivating Example

Consider a web search provider, which want to present a result list of web pages without duplicates. To achieve that, for every pair of web pages, we drop one of them, if they are textually very similar. This is the case when their *resemblance* R is greater than some threshold R_0 . But to be able to use the *resemblance* as a measurement, we have to map each page to a set. One could imagine to define this mapping from the page to the set of all words occurring on the page. As in several studies [3, 4] we will instead map a page to a set of *shingles*. A shingle is a string of w contiguous words, and we include a *shingle* in our result set if the *shingle* occurs on the page (in the same order). Typically we choose $w = 5$. Figure 2 shows an exemplary mapping.



Figure 2: The mapping from a web page to a set of *shingles*.

Clearly, the number of possible shingles and therefore D is huge. Assume 10^5 different English words, we have $D = (10^5)^5$. Thus computing the exact simi-

larities for all pairs of pages of a web search would require prohibitive storage. For a single pair, we would already need $\Theta(D)$ storage.

1.3 Original Minwise Hashing

A working approximative solution to this problem was described by Broder and his colleagues [3, 4]. Their algorithm is based on the following observation:

Suppose a random permutation π is performed on Ω^1 ,

$$\pi : \Omega \longrightarrow \Omega, \quad \text{where } \Omega = \{0, 1, \dots, D-1\}.$$

To simplify notation, we overload the definition of π to work also for subsets of Ω . Thus with $\pi(S_n)$ we denote the application of the permutation π to every element of the set S_n . More precisely,

$$\begin{aligned} \pi : 2^\Omega &\longrightarrow 2^\Omega, \\ \pi : S_n &\longmapsto \pi(S_n) = \{\pi(i) | i \in S_n\}. \end{aligned}$$

A key role play the minima of the permuted sets:

$$e_{S_n, \pi} = \min(\pi(S_n))$$

We use $\Pr[\cdot]$ as a shorthand for the *Probability* of some event. An elementary probability argument shows that for two sets $X, Y \subseteq \Omega$

$$\Pr[e_{X, \pi} = e_{Y, \pi}] = \Pr[\min(\pi(X)) = \min(\pi(Y))] = \frac{|X \cap Y|}{|X \cup Y|} = R. \quad (2)$$

We can now build an unbiased estimator \hat{R}_M of R with k minwise independent permutations, $\pi_1, \pi_2, \dots, \pi_k$:

$$\hat{R}_M = \frac{1}{k} \sum_{j=1}^k 1 \{e_{X, \pi_j} = e_{Y, \pi_j}\}, \quad (3)$$

$$\text{Var}(\hat{R}_M) = \frac{1}{k} R(1 - R). \quad (4)$$

We later refer to e_{S_n, π_j} as a **sample** and to k as the **sample size**. To store a **sample** may require e.g., 64 bits. Note that the minimum function $\min(\cdot)$ needs $O(D)$ time and constant space and can be precomputed for each set S_n and permutation π_j .²

¹We assume there is a perfect hash function applied to the elements of the original domain which always gives us $\Omega = \{0, 1, \dots, D-1\}$. Note that in the paper [1], the hash function is applied after the permutation and the minimum-function, but it is simpler to understand this way.

²A detailed runtime analysis will follow for the improved algorithm in section 2.4 on page 6

Many applications, especially duplicate detection, are interested in detecting somewhat high similarity, thus an approximation seems reasonable. From equation (4) on the previous page we learn, that the accuracy can be adjusted by choosing the **sample size** appropriately.

However, finding duplicates out of m objects, e.g. web pages, still needs $O(m^2)$ comparisons. There are number of approaches to deal with that problem, but we will not discuss it here and refer to [4, 5].

2 b-bit k-permutation Minwise Hashing

We will present an algorithm which improves on the storage requirements of the original MINWISE HASHING. The idea is to reduce the size of each **sample** by only taking b bits of it, as opposed to, e.g. 64 bits. Intuitively, this will increase the estimation variance $Var(\hat{R}_M)$, at the same **sample size** k . To maintain the same accuracy, we have to increase k . One can show, if the resemblance is not too small, we will not have to increase k much and in total use less storage.

For example, when $b = 1$ and $R = 0.5$, the estimation variance will increase at most by a factor of 3. To keep the same accuracy, we have to increase the **sample size** by a factor of 3. If we before stored each minimum $\min(\pi_j(S_n))$ using 64 bits, the improvement with $b = 1$ is $64/3 = 21.3$.

Consider again two sets $X, Y \subseteq \Omega$, on which a random permutation $\pi : \Omega \rightarrow \Omega$ is applied. We extend our notation with

$$e_{X,b,\pi} = b \text{ lowest bits of } \min(\pi(X))$$

$$e_{Y,b,\pi} = b \text{ lowest bits of } \min(\pi(Y))$$

Theorem 1. *Assume D is large.*

$$P_b = \Pr[1 \{e_{X,b,\pi} = e_{Y,b,\pi}\}] = C_{1,b} + (1 - C_{2,b})R \quad (5)$$

$$r_X = \frac{|X|}{D}, \quad r_Y = \frac{|Y|}{D} \quad (6)$$

$$C_{1,b} = A_{X,b} \frac{r_Y}{r_X + r_Y} + A_{Y,b} \frac{r_X}{r_X + r_Y}, \quad (7)$$

$$C_{2,b} = A_{X,b} \frac{r_X}{r_X + r_Y} + A_{Y,b} \frac{r_Y}{r_X + r_Y}$$

$$A_{X,b} = \frac{r_X [1 - r_X]^{2^b - 1}}{1 - [1 - r_X]^{2^b}}, \quad A_{Y,b} = \frac{r_Y [1 - r_Y]^{2^b - 1}}{1 - [1 - r_Y]^{2^b}} \quad (8)$$

The intuition for the additional terms $C_{1,b}$ and $C_{2,b}$ in (5) compared to (2) on the previous page is that we have to account for a type of “false positive”: When two minima agree on their last b bits, $e_{X,b,\pi} = e_{Y,b,\pi}$, it’s still possible that their are different, $e_{X,\pi} \neq e_{Y,\pi}$. Thus even when $R = 0$, the collision

probability P_b is not zero, but rather $C_{1,b}$. This makes the derivation much more complicated.³

Even though D is assumed to be large, experiments show that even for $D = 20$ the absolute error caused by using (5) on the previous page is < 0.01 .

2.1 The Estimator

From (5) on the preceding page of Theorem 1 on the previous page we derive the estimator \hat{R}_b for R :

$$\hat{R}_b = \frac{\hat{P}_b - C_{1,b}}{1 - C_{2,b}} \quad (9)$$

$$\hat{P}_b = \frac{1}{k} \sum_{j=1}^k 1 \{e_{X,b,\pi_j} = e_{Y,b,\pi_j}\} \quad (10)$$

This estimator is unbiased, i.e. $E[\hat{R}_b] = R$. Furthermore, the variance of \hat{R}_M converges to the variance of \hat{R}_b , i.e.

$$\lim_{b \rightarrow \inf} \text{Var}(\hat{R}_b) = \frac{R(1-R)}{k} = \text{Var}(\hat{R}_M) \quad (11)$$

2.2 The Algorithm

Based on the theoretical results, Algorithm 1 presents the procedure of b -bit (k -permutation) Hashing.

Algorithm 1 B-BIT MINWISE HASHING algorithm, applied to estimating pairwise resemblances in a collection of N sets.

Input: Sets $S_n \subseteq \Omega = \{0, 1, \dots, D-1\}$, $n = 1$ to N . $\triangleright D = |\Omega|$

Output: Estimated resemblance \hat{R}_b

// Pre-processing

Generate k random permutations $\pi_j : \Omega \rightarrow \Omega$, $j = 1$ to k

for all $n = 1$ to N , $j = 1$ to k **do**

Store the lowest b bits of $\min(\pi_j(S_n))$, denoted by e_{S_n,b,π_j} .

end for

// Estimation (Use two sets X, Y as an example)

Compute $\hat{P}_b = \frac{1}{k} \sum_{j=1}^k 1 \{e_{X,b,\pi_j} = e_{Y,b,\pi_j}\}$

Estimate the resemblance by $\hat{R}_b = \frac{\hat{P}_b - C_{1,b}}{1 - C_{2,b}}$, where $C_{1,b}$ and $C_{2,b}$ are from Theorem 1 on the previous page

³A proof of Theorem 1 on the preceding page can be found in the appendix of [6].

2.3 Deriving the Hamming Distance

Another well-known measurement for the similarity is the *hamming distance*. For the purpose of calculating the *hamming distance* between two sets $X, Y \subseteq \Omega = \{0, 1, \dots, D-1\}$, the sets are first mapped to a D -dimensional binary vector x, y resp.:

Definition 2. Let vector $x, y \in \{0, 1\}^D$, $x_i = 1 \{i \in X\}$, $y_i = 1 \{i \in Y\}$. The hamming distance between X and Y is

$$H = \sum_{i=0}^{D-1} [x_i \neq y_i] = |X \cup Y| - |X \cap Y| = |X| + |Y| - 2a \quad (12)$$

If we reformulate 1 on page 2 as

$$a = \frac{R}{1+R}(|X| + |Y|), \quad (13)$$

we can use the B-BIT MINWISE HASHING algorithm to estimate H with

$$\hat{H}_b = |X| + |Y| - 2 \frac{\hat{R}_b}{1 + \hat{R}_b}(|X| + |Y|) = \frac{1 - \hat{R}_b}{1 + \hat{R}_b}(|X| + |Y|) \quad (14)$$

Experiments show that this approach is significantly faster than standard methods for computing the *hamming distance*.

2.4 Drawbacks

The major problem of B-BIT MINWISE HASHING is the costly preprocessing. Consider an application in machine learning, where the sets S_n represents some properties of an object (e.g. a document or an image). Often one wants to add objects dynamically at runtime and compare them to other objects by finding the *resemblance* of their properties. Finding the k minima under the permutations may take too long, i.e. uses $O(k|S_n|)$ to $O(kD)$ time per set. In general, the costly preprocessing may cause problems in user-facing applications, where the testing efficiency for new data objects is crucial. There is the need for a entirely fast algorithm to keep this applications responsive.

Another drawback is that storing k permutations is sometimes impractical. If e.g. $D = 10^9$, one permutation vector uses 4GB, which is still possible to store. But the space needed to store e.g. $k = 500$ permutation vectors (each of length D), is not tolerable.

3 One Permutation Hashing

This algorithm is directly motivated by the optimization potential of the standard MINWISE HASHING method: intuitively, it ought to be “wasteful” in that

all elements in a set are permuted, scanned but only the minimum will be used. As the name already suggests, we reduce the preprocessing step to only one permutation.

As in section 2.3 on the previous page we will represent sets $S_n \subseteq \Omega$ as vectors $s_n \in \{0, 1\}^D$, $(s_n)_i = 1 \{i \in S_n\}$. We will setup a running example with $X, Y, Z \subseteq \Omega = \{0, 1, \dots, 15\}$. Let

$$X = \{2, 4, 7, 13\}, \quad Y = \{0, 3, 6, 13\}, \quad Z = \{0, 1, 10, 12\}.$$

Now we build up a data matrix where the rows are equal to the vector representations of the permuted sets:

$$\begin{array}{l} \text{Sets:} \end{array} \quad \overbrace{\begin{array}{c} 0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11 \ 12 \ 13 \ 14 \ 15 \\ \pi(X): \begin{pmatrix} 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \\ \pi(Y): \begin{pmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \\ \pi(Z): \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \end{pmatrix} \end{array}}^{\Omega} \quad (15)$$

The idea is to divide the columns evenly into t (here $t = 4$) bins (parts), take the minimum in each bin. Because later we only compare minima within one bin, we can re-index the elements to use the smallest possible representation:

$$\begin{array}{l} \text{Sets:} \end{array} \quad \overbrace{\begin{array}{c} \Omega' \quad \Omega' \quad \Omega' \quad \Omega' \\ 0 \ 1 \ 2 \ 3 \ 0 \ 1 \ 2 \ 3 \ 0 \ 1 \ 2 \ 3 \ 0 \ 1 \ 2 \ 3 \\ \pi(X): \begin{pmatrix} 0 & 0 & \textcircled{1} & 0 & \textcircled{1} & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & \textcircled{1} & 0 & 0 \\ \textcircled{1} & 0 & 0 & 1 & 0 & 0 & \textcircled{1} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcircled{1} & 0 & 0 \\ \textcircled{1} & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \textcircled{1} & 0 & \textcircled{1} & 0 & 0 & 0 & 0 \end{pmatrix} \\ \pi(Y): \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \\ \pi(Z): \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \end{array}}^{\Omega'} \quad (16)$$

bin 1
bin 2
bin 3
bin 4

We get the minima-vectors

$$\begin{aligned} v_X &= [2, 0, *, 1], \\ v_Y &= [0, 2, *, 1], \\ v_Z &= [0, *, 2, 0], \end{aligned} \quad (17)$$

where '*' denotes an empty bin.

To derive the *resemblance* between two sets, e.g. X, Y , we introduce two definitions:

$$\begin{aligned} \text{number of "jointly empty bins":} \quad N_{emp} &= \sum_{j=1}^t I_{emp,j}, \\ \text{number of "matched bins":} \quad N_{mat} &= \sum_{j=1}^t I_{mat,j}, \end{aligned} \quad (18)$$

where $I_{emp,j}$ and $I_{mat,j}$ are defined for the j -th bin, as

$$\begin{aligned} I_{emp,j} &= \begin{cases} 1 & \text{if both } \pi(X) \text{ and } \pi(Y) \text{ are empty in the } j\text{-th bin} \\ 0 & \text{otherwise} \end{cases} \\ I_{mat,j} &= \begin{cases} 1 & \text{if both } \pi(X) \text{ and } \pi(Y) \text{ are not empty and the smallest} \\ & \text{elements in the } j\text{-th bin matches, i.e. } (v_X)_j = (v_Y)_j \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (19)$$

3.1 The Estimator

Recall the notations $a = |X \cap Y|$ and $f = |X \cup Y| = |X| + |Y| - a$. We formulate the estimator as

Lemma 1. *The resemblance is estimated by*

$$\hat{R}_{mat} = \frac{N_{mat}}{k - N_{emp}} \quad (20)$$

is unbiased, i.e.

$$\mathbb{E}[\hat{R}_{mat}] = R. \quad (21)$$

In our example we have $N_{emp} = 1$ and $N_{mat} = 1$. Thus $\hat{R}_{mat} = 1/3$.

Because it is a bit surprising that the estimator is unbiased, we give a

Proof of Lemma 1. Because we assume that the sets and thus the data vectors are not completely empty, it holds

$$t - N_{emp} > 0 \Rightarrow P[t - N_{emp} > 0] = 1,$$

hence we get rid of division-by-zero problems and $m > 0$ in (24) is always true. From the definitions in (19) follows

$$I_{emp,j} = 1 \Rightarrow I_{mat,j} = 0, \quad (22)$$

$$\mathbb{E}[I_{mat,j} | I_{emp,j} = 0] = R, \quad (23)$$

$$\text{for } m > 0: \quad \mathbb{E}[I_{mat,j} | t - N_{emp} = m] = \frac{m}{t} R, \quad (24)$$

and we derive

$$\begin{aligned} \mathbb{E}[N_{mat} | t - N_{emp} = m] &\stackrel{(22)}{=} \sum_{j=1}^{t-N_{emp}} \mathbb{E}[I_{mat,j} | t - N_{emp} = m] \\ &\stackrel{(23)}{=} \sum_{j=1}^{t-N_{emp}} R = R(t - N_{emp}). \end{aligned} \quad (25)$$

Finally,

$$\begin{aligned} \mathbb{E} \left[\frac{N_{mat}}{t - N_{emp}} \middle| t - N_{emp} = m \right] &\stackrel{(25)}{=} R, \quad (\text{independent of } N_{emp}) \\ \Rightarrow \mathbb{E} \left[\frac{N_{mat}}{t - N_{emp}} \right] &= \mathbb{E} [\hat{R}_{mat}] = R \end{aligned} \quad (26)$$

□

We give the variance without a proof.⁴

$$\text{Var} [\hat{R}_{mat}] = R(1 - R) \left(\mathbb{E} \left[\frac{1}{t - N_{emp}} \right] \left(1 + \frac{1}{f - 1} \right) - \frac{1}{f - 1} \right) \quad (27)$$

If we have very few empty bins, i.e. N_{emp} is essentially zero and $t \ll f$, the term simplifies to

$$\text{Var} [\hat{R}_{mat}] \approx \frac{R(1 - R)}{t} \left(\frac{f - t}{f - 1} \right) \approx \frac{R(1 - R)}{t} \quad (28)$$

as expected.

3.2 The Algorithm

Based on the theoretical results, Algorithm 2 presents the procedure of ONE PERMUTATION HASHING.

Algorithm 2 ONE PERMUTATION HASHING algorithm, applied to estimating pairwise resemblances in a collection of N sets.

Input: Sets $S_n \subseteq \Omega = \{0, 1, \dots, D - 1\}$, $n = 1$ to N . ▷ $D = |\Omega|$

Output: Estimated resemblance \hat{R}_{mat}

// Pre-processing

Generate one random permutations $\pi : \Omega \rightarrow \Omega$

for all $n = 1$ to N **do**

Permute set S_n with π , store the re-indexed minima of each of the t bins in a data vector $v_{S_n} \in \{0, \dots, \lfloor |S_n|/t \rfloor\}^t$. ▷ see (16)

end for

// Estimation (Use two sets X, Y as an example)

Estimate the resemblance by $\hat{R}_{mat} = \frac{N_{mat}}{t - N_{emp}}$ ▷ see (18)

Empirical results show that the ONE PERMUTATION HASHING scheme performs as well or even slightly better than the B-BIT K-PERMUTATION HASHING scheme.

Let's have a look at the runtime and storage requirements and compare them with Algorithm 1 on page 5. The pre-processing step generates one permutation vector, which uses $O(D \log(D))$ storage, which even for large D is still

⁴A proof can be found in the appendix of [2].

practical (in contrast to storing k permutation vectors for the MINWISE HASHING method). The pre-processing also generates a data vector with t elements (which uses themselves only $O(\log(D/t))$ bits because we re-indexed them within the bins) for each set, thus $O(t \log(D/t))$ bits per set.

Table 1 summarizes the runtime and storage requirements of the different algorithms.

Space used for ...	storing permutation(s)	storing pre-processing data per set
Naïve approach (store whole sets)	—	$O(D)$
original MINWISE HASHING	$O(kD \log(D))$	$O(k \log(D))$
B-BIT MINWISE HASHING	$O(kD \log(D))$	$O(kb)$
ONE PERMUTATION HASHING	$O(D \log(D))$	$O(t \log(D/t))$

Table 1: Algorithm comparison. We only list the space complexity because the processing time boundaries are given by the time needed to read the stored data at least once and hence is dominated by the space complexity.

4 Conclusions

MINWISE HASHING and ONE PERMUTATION HASHING are standard techniques for efficiently estimating set similarity in massive datasets. We started at explaining the original MINWISE HASHING and demonstrated how reducing the amount of stored bits per **sample** to b bits leads to a effective reduction of the required storage and consequently computational overhead.

We then went over to ONE PERMUTATION HASHING, which uses the same basic idea, but makes better use of the permuted set and therefore only needs one permutation. Overall, this last algorithm is superior to the others in many aspects.

References

- [1] Ping Li 0001 and Arnd Christian König. Theory and applications of b -bit minwise hashing. *Commun. ACM*, 54(8):101–109, 2011.
- [2] Ping Li 0001, Art B. Owen, and Cun-Hui Zhang. One permutation hashing for efficient search and learning. *CoRR*, abs/1208.1259, 2012.
- [3] Andrei Z. Broder. On the resemblance and containment of documents. In *Compression and Complexity of Sequences, IEEE Computer Society*, pages 21–29, Italy, 1997. Positano.

- [4] Andrei Z. Broder, Steven C. Glassman, Mark S. Manasse, and Geoffrey Zweig. Syntactic clustering of the web. *Computer Networks*, pages 1157–1166, 1997.
- [5] Moses Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the 34th Annual ACM Symposium on Theory of Computing (STOC-02)*, pages 380–388, New York, May 19–21 2002. ACM Press.
- [6] Ping Li 0001 and Arnd Christian König. b-bit minwise hashing. *CoRR*, abs/0910.3349, 2009.