

1.2 Motivation

Although MinHash has many variants, the fundamental idea behind it is to construct an event, in that case, the collision of minimum hash values in S_1, S_2 , whose probability of happening equals $J_k(S_1, S_2)$, the quantity MinHash tries to estimate. Once we have this event, we can estimate the probability of it happening by repeating experiments, thus obtaining an estimate of $J_k(S_1, S_2)$. We will try to follow the same idea to come up with a sketching algorithm that estimates the Weighted Jaccard Similarity. We know that

$$\begin{aligned} J_k^w(S_1, S_2) &= \frac{\sum_{\sigma} \min\{N_k(\sigma, S_1), N_k(\sigma, S_2)\}}{\sum_{\sigma} \max\{N_k(\sigma, S_1), N_k(\sigma, S_2)\}} \\ &= \frac{|T_k(S_1) \cap T_k(S_2)|}{|T_k(S_1) \cup T_k(S_2)|} \end{aligned}$$

So we need an event whose probability of happening equals $\frac{|T_k(S_1) \cap T_k(S_2)|}{|T_k(S_1) \cup T_k(S_2)|}$. Of course, we also need the event to be easy to “simulate” so that running the experiments repeatedly can be done using little resources. An obvious event whose probability of happening equals $J_k^w(S_1, S_2)$ is for a uniform random sample of indexed k -mer (σ, i) in $T_k(S_1) \cup T_k(S_2)$ to fall in $T_k(S_1) \cap T_k(S_2)$. The problem is that in order to perform such sampling, we need access to the indexed k -mer sets $T_k(S_1)$ and $T_k(S_2)$, which may not always be available. More specifically, in order to sample from $T_k(S_1) \cup T_k(S_2)$, we need the input strings S_1, S_2 to contain information about all the k -mer indexes. Pictorially, we need S_1, S_2 to have the following forms:

$$\begin{aligned} S_1 &= \dots (\sigma, 1) \dots (\sigma, 2) \dots (\sigma, l) \dots (\sigma, m) \dots \\ S_2 &= \dots (\sigma, 1) \dots (\sigma, 2) \dots (\sigma, l) \dots \end{aligned}$$

which would allow us to somehow “pair up” the first $l = \min\{N_k(\sigma, S_1), N_k(\sigma, S_2)\}$ occurrences of σ one by one and treat each pair as only “one occurrence” to obtain $T_k(S_1) \cap T_k(S_2)$. In other words, *certain pairs* σ should not be distinguished.

But in reality, the input strings may well not be in this form and simply appear as sequences of letters:

$$\begin{aligned} S_1 &= \dots \sigma \dots \sigma \dots \sigma \dots \sigma \dots \\ S_2 &= \dots \sigma \dots \sigma \dots \sigma \dots \end{aligned}$$

Clearly, we can distinguish all occurrences of σ in S_1 from those in S_2 , and distinguish all occurrences of σ in S_1 by their (absolute) position in the string. But this would treat all occurrences of σ as distinct, and we would have no way of knowing which pairs of σ in S_1, S_2 should be treated as the same. Note that going through the strings to compute all the occurrence numbers would be computationally infeasible since we assumed the k -mer sets are huge and input strings massive.

Although we do not directly have access¹ to $T_k(S_1)$ and $T_k(S_2)$, by treating all k -mer occurrences as distinct, we can still “sample” an indexed k -mer in the following way:

1. Randomly sample a position from all the available positions in S_1 or S_2 where a k -mer might reside.
2. If σ is the k -mer at that position, we can go through the string to compute the occurrence number of m of this particular σ by counting how many times σ appears in the string up to its position.

Moreover, once (σ, m) is obtained, and the random position from the first step is in, say, S_1 , then it directly follows from the definition of indexed k -mers that we can easily check whether (σ, m) is in $T_k(S_1) \cap T_k(S_2)$ by going through S_2 and count whether σ appears at least m times in S_2 .

This method provides a randomly sampled indexed k -mer (σ, m) , and an easy way of checking whether it is in $T_k(S_1) \cap T_k(S_2)$. The problem is that we did not directly sample (σ, m) from $T_k(S_1) \cup T_k(S_2)$. Remember that (σ, m) is determined by randomly sampling from all available positions where a k -mer might occur. In other words, the true sample space is the set of all k -mer occurrences treated as distinct, which has size $\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)$, as opposed to $T_k(S_1) \cup T_k(S_2)$, which has size $\sum_{\tau} \max\{N_k(\tau, S_1), N_k(\tau, S_2)\}$. So the actual sampling process is not a uniform sampling from $T_k(S_1) \cup T_k(S_2)$, which means the probability that $(\sigma, i) \in T_k(S_1) \cap T_k(S_2)$ is likely not $J_k^w(S_1, S_2)$.

But there is a silver line: we know exactly how larger the true sample space is compared to $T_k(S_1) \cup T_k(S_2)$. So we might hope that we can somehow use this information to transform the probability $\mathbb{P}[(\sigma, i) \in T_k(S_1) \cap T_k(S_2)]$ into $J_k^w(S_1, S_2)$. If this is true, then as long as the process of transformation does not introduce too much error, we can still indirectly estimate $J_k^w(S_1, S_2)$ by estimating the aforementioned probability. As we will see below, this hope is indeed true.

1.3 The Algorithm

Now we present the algorithm based on the idea discussed in the previous section and confirm that it indeed works.

We now show that the algorithm indeed estimates $J_k^w(S_1, S_2)$. As previously discussed, f is the indicator for the event that the random sample (σ, m) falls in $T_k(S_1) \cap T_k(S_2)$. Because (σ, m) is not directly sampled from $T_k(S_1) \cup T_k(S_2)$, $\mathbb{P}[f \text{ returns } 1]$ will not be $J_k^w(S_1, S_2)$. But we can compute $\mathbb{P}[f \text{ returns } 1]$ explicitly as follows: fix any k -mer σ , and without loss of generality, suppose S_1, S_2 have the following form:

$$\begin{aligned} S_1 &= \sigma \dots \sigma \dots \sigma \dots \textcolor{blue}{\sigma} \dots \textcolor{red}{\sigma} \dots \\ S_2 &= \sigma \dots \sigma \dots \sigma \dots \textcolor{blue}{\sigma} \dots \end{aligned}$$

¹However, if for some reason we do have access to $T_k(S_1), T_k(S_2)$, we can actually estimate an even better similarity measurement. This will be discussed in Section ??.

Algorithm 1 $f(S_1, S_2, k)$: experiment of random sampling

1. Sample a position i uniformly at random from all possible positions in S_1, S_2 where a k -mer resides.
 2. Let $S :=$ the string containing position i , $S' :=$ the other string. $\sigma :=$ the k -mer at position i in S .
 3. Pass through S , find $m \leftarrow$ number of occurrences of σ in S up to and including position i .
 4. Pass through S' , check if σ appears at least m times in S' . Return 1 if true, 0 otherwise.
-

Algorithm 2 $g(S_1, S_2, k)$: repeat experiments f to estimate $J_k^w(S_1, S_2)$

1. Run R independent copies of $f(S_1, S_2, k)$. $\delta_j =$ return value of the j th copy.
 2. $\hat{p} = \frac{1}{R} \sum_{j=1}^R \delta_j$.
 3. Output $\frac{\hat{p}}{2-\hat{p}}$ as the estimate of $J_k^w(S_1, S_2)$.
-

where σ and σ mark the last occurrence of σ in S_1, S_2 , respectively. Then by definition, σ is the $N_k(\sigma, S_2)$ th copy of σ in S_2 . Let σ be the $N_k(\sigma, S_2)$ th copy of σ in S_1 .

Clearly, the event f returns 1 and the sampled k -mer in step 2 of f equals σ happens if and only if the sample is a σ that falls in the following region:

$$\begin{array}{l} S_1 = \sigma \dots \sigma \dots \sigma \dots \sigma \dots \sigma \dots \sigma \dots \\ S_2 = \sigma \dots \sigma \dots \sigma \dots \sigma \dots \sigma \dots \sigma \dots \end{array}$$

The circled region above contains $2 \cdot N_k(\sigma, S_2) = 2 \min\{N_k(\sigma, S_1), N_k(\sigma, S_2)\}$ occurrences of σ . As argued before, the true sample space is equivalently all k -mer occurrences treated distinct and thus has size $\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)$. Because the sampling in step 1 of f is uniform sampling, this means

$$\mathbb{P}[f \text{ returns 1 and sample is } \sigma] = \frac{2 \min\{N_k(\sigma, S_1), N_k(\sigma, S_2)\}}{\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)}$$

By marginalizing out σ , we get:

$$\begin{aligned}
p &= \mathbb{P}[f \text{ returns } 1] \\
&= \sum_{\sigma} \mathbb{P}[f \text{ returns } 1 \text{ and sample is } \sigma] \\
&= \sum_{\sigma} \frac{2 \min\{N_k(\sigma, S_1), N_k(\sigma, S_2)\}}{\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)} \\
&= \frac{2 \sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\}}{\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)}
\end{aligned}$$

In comparison:

$$J_k^w(S_1, S_2) = \frac{\sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\}}{\sum_{\tau} \max\{N_k(\tau, S_1), N_k(\tau, S_2)\}}$$

which is different from p , as expected. But we can easily transform:

$$\begin{aligned}
p &= \frac{2 \sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\}}{\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)} \\
&= \frac{2 \sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\}}{\sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\} + \max\{N_k(\tau, S_1), N_k(\tau, S_2)\}}
\end{aligned}$$

and solve

$$\begin{aligned}
\frac{p}{2-p} &= \frac{\sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\}}{\sum_{\tau} \max\{N_k(\tau, S_1), N_k(\tau, S_2)\}} \\
&= J_k^w(S_1, S_2)
\end{aligned}$$

thus justifying the estimator $\hat{p} = \frac{1}{R} \sum_{j=1}^R \delta_j$ in $g(S_1, S_2, k)$.

1.4 Accuracy Analysis

Let us analyze the accuracy of the estimation. First note that

$$\mathbb{E}[\hat{p}] = \frac{1}{R} \sum_{j=1}^R \mathbb{E}[\delta_j] = p$$

So \hat{p} is an unbiased estimator of the true probability p that f returns 1. We proved that p satisfies

$$\frac{p}{2-p} = J_k^w(S_1, S_2)$$

So we can hope $\frac{\hat{p}}{2-\hat{p}}$ will not be too far from $J_k^w(S_1, S_2)$ if \hat{p} is not too far from p .
Indeed:

$$\begin{aligned} \left| \frac{\hat{p}}{2-\hat{p}} - J_k^w(S_1, S_2) \right| &= \left| \frac{p}{2-p} - \frac{\hat{p}}{2-\hat{p}} \right| \\ &= \frac{2|p-\hat{p}|}{(2-p)(2-\hat{p})} \\ &\leq 2|p-\hat{p}| \\ &\quad (\text{Since } 1 < 2-p, 2-\hat{p} < 2.) \end{aligned}$$

Thus if \hat{p} is within ϵ of p , then $\frac{\hat{p}}{2-\hat{p}}$ will be within 2ϵ of $J_k^w(S_1, S_2)$. Recall from the introduction the goal we want to achieve with respect to accuracy: we want the probability that the estimator $\frac{\hat{p}}{2-\hat{p}}$ is far from $J_k^w(S_1, S_2)$ to be small. We can easily formalize this notion as follows. First, compute

$$\begin{aligned} \text{Var}[\hat{p}] &= \text{Var}\left[\frac{1}{R} \sum_{j=1}^R \delta_j\right] \\ &= \frac{1}{R^2} \sum_{j=1}^R \text{Var}[\delta_j] \quad (\text{by independence}) \\ &= \frac{1}{R} (\mathbb{E}[\delta_1^2] - \mathbb{E}[\delta_1]^2) \\ &= \frac{p - p^2}{R} \end{aligned}$$

Then by Chebyshev's inequality, for $\epsilon > 0$:

$$\begin{aligned} \mathbb{P}[|\hat{p} - \mathbb{E}[\hat{p}]| \geq \epsilon] &= \mathbb{P}[|\hat{p} - p| \geq \epsilon] \\ &\leq \frac{\text{Var}[\hat{p}]}{\epsilon^2} \\ &= \frac{p - p^2}{R\epsilon^2} \end{aligned}$$

Combine this with $\left| \frac{\hat{p}}{2-\hat{p}} - J_k^w(S_1, S_2) \right| \leq 2|p-\hat{p}|$, we get:

$$\begin{aligned} \mathbb{P}\left[\left| \frac{\hat{p}}{2-\hat{p}} - J_k^w(S_1, S_2) \right| < \epsilon\right] &\geq \mathbb{P}[|\hat{p} - p| < \epsilon/2] \\ &= 1 - \mathbb{P}[|\hat{p} - p| \geq \epsilon/2] \\ &\geq 1 - \frac{4(p - p^2)}{R\epsilon^2} \end{aligned}$$

From the above bound, it is clear that given a fixed ϵ , we can make the probability that our estimate falls within ϵ of the true $J_k^w(S_1, S_2)$ by running the “experiment” $f(S_1, S_2, k)$ enough times, which conforms with our intuition from Section ??.

1.5 Resource analysis and streaming model considerations

Recall from Section ?? that apart from the accuracy requirement, we also need to make sure that the algorithm does not take too much computing resources, primarily time and space. Since g simply runs R independent copies of f , the resources it takes should be approximately R times that of f , or possibly the same as f , depending on whether we run the copies of f in parallel. We can divide f into 3 major steps:

1. Sample a k -mer σ at a random position i .
2. Pass through one of S_1, S_2 once to compute the occurrence number m of σ .
3. Pass through another string once to check if $(\sigma, m) \in T_k(S_1) \cap T_k(S_2)$.

Along each pass, we need to constantly compare the sample k -mer with another k -mer in the string. Since comparing two k -mers takes $O(k)$ time and the total number of k -mers along the way is $O(|S_1| + |S_2|)$, f takes $O(k(|S_1| + |S_2|))$ time to execute steps 2) and 3). If we know the length of the two strings $|S_1|, |S_2|$, then sampling a random position is equivalent to simply generating a random integer i in the range $[0, |S_1| - k] \cup [|S_1|, |S_1| + |S_2| - k]$, which can be easily done using $O(1)$ time and space. However, the algorithm needs to retrieve the k -mer at position i . And the time it does this is dependent upon how the input strings S_1, S_2 are stored.

First, there is the perhaps trivial situation when S_1, S_2 are already stored in Random Access Memory, which means we can access the letter at any position in a string and retrieve the k -mer at position i in $O(1)$ time. In this situation, it is often true that information about $|S_1|, |S_2|$ is easily accessible (for example, Strings in Python keep track of their own lengths, which can be found in $O(1)$ time). In this situation, f only needs to pass through S_1, S_2 once and the runtime would be $O(k(|S_1| + |S_2|))$. If $|S_1|, |S_2|$

However, this situation assumes that we must already have enough memory to store S_1, S_2 , which seems to defeat the point of using probabilistic sketching. What about the more common streaming model, where the input strings appear as little pieces in an incoming stream?

Without loss of generality, we can assume the strings appear as streams of k -mers because otherwise we can simply wait until enough letters have arrived to assemble a k -mer. In this model, we need to perform uniform sampling of a k -mer occurrence from an incoming stream with an unknown number of items.

Of course, there is one easy and trivial solution:

1. Maintain 2 counters and go through the stream once to find out $|S_1|, |S_2|$.
2. Sample a random integer as the position i , go through the stream again to obtain the k -mer at position i .
3. Go through the stream for a third time to do part(2) of f .

This method does not increase the asymptotic runtime of f (still $O(k(|S_1| + |S_2|))$). But ideally we want to restart the stream as few times as possible because it might be

expensive in practical situations (for example, sending large amount of data over the internet), and worse still, in some cases we only have access to the stream once!

We can use a common technique called Reservoir Sampling to reduce the number of passes required for sampling to one. In general, Reservoir Sampling can sample without replacement l items from a stream of unknown length n such that each one of the n items has probability l/n of ending up in the final sample. Since f only requires sampling one position with probability $1/n$, we can use the most basic version of Reservoir Sampling, which goes as follows:

- Start the stream and keep the first item.
- As the i th item in the stream arrives, update the sample to be the i th item with probability $1/i$ and do nothing with probability $1 - 1/i$.

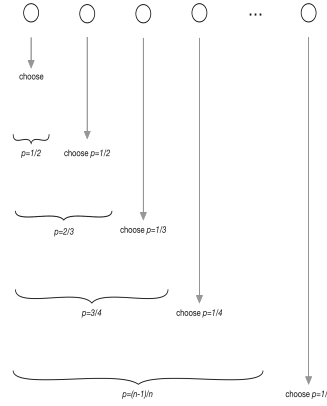


Figure 1: The process of Reservoir Sampling.

It is easy to see that

$$\forall t, \mathbb{P}[\text{select item } t] = \frac{1}{t} \cdot \frac{t}{t+1} \cdot \frac{t+1}{t+2} \cdots \frac{n-1}{n} = \frac{1}{n}$$

So with Reservoir Sampling, the experiment function f can be implemented as:

Despite the improvement, this implementation still requires two passes through the stream, which is less than ideal and greatly limits the applicability of the algorithm. The author has been attempting to try to reduce the number of passes needed to 1, a discussion of which will be provided in the next section.

Finally, it is clear that f only needs to store the sample k -mer σ , another k -mer in the stream to compare with σ , and a constant number of counters. So f requires only $O(k)$ space.

1.6 An ongoing attempt to further improve f

Let us explore some possible directions to improve f .

Algorithm 3 $f+(S_1, S_2, k)$: Improvements via Reservoir Sampling.

1. Use Reservoir Sampling to sample a k -mer σ (and its position) uniformly at random.
 2. Restart the stream again. Count the occurrence number m of σ in the string from which it is sampled.
 3. Count whether the other string contains at least m copies of σ , and return 0 or 1 accordingly.
-

Recall the three major steps of f :

1. Sample a k -mer σ at a random position i .
2. Pass through one of S_1, S_2 once to compute the occurrence number m of σ .
3. Pass through another string once to check if $(\sigma, m) \in T_k(S_1) \cap T_k(S_2)$.

The difficulty of reducing the required number of passes through the stream by f arises from the fact that in the incoming stream model, we need one pass to execute step 1) (i.e., obtaining the sketch σ and i), and another pass to finish steps 2) and 3) (i.e., utilizing the sketch). In comparison, once MinHash obtains its sketch, namely the two minimum hash values, we no longer need information regarding the input strings to use the sketch (the two minimum hash values can be compared in $O(1)$ time).

Obviously, in order to reduce the number of passes, we need to somehow perform one of the steps above without passing through the entire input strings. For example, f only requires one pass when S_1, S_2 are both stored in RAM, because in that case step 1) reduces to generating random index and string indexing. This is of course unfeasible in the incoming stream model, so we are motivated to consider steps 2) and 3).

One possible idea is the following. Since the occurrence number m of an indexed k -mer (σ, m) is simply an integer within the range $[1, N_k(\sigma, S)]$, one might think that if σ and its frequency $N_k(\sigma, S)$ are given, then a random sample of an indexed k -mer can be obtained by randomly generating an integer in the range $[1, N_k(\sigma)]$. Of course, we might be given all the k -mer frequencies, and we assumed that finding them by brute force is not feasible. But there is a sketching algorithm called Count-min sketch [?] that can serve as a “frequency table” with limited space and time, and of course, some chances of errors.

Without getting into the details, Count-min sketch consumes a stream of items, and at any time, we can make a query to find out the approximate frequency of any item σ it has seen so far. The approximate frequency might be larger than the true frequency of that item, but it will not be smaller. One can tune the parameters of Count-min sketch to make the probability of the approximate frequency being too far from the true frequency low. If we combine Count-min sketch with the Reservoir sampling technique, we obtain the following algorithm:

Algorithm 4 $f_{++}(S_1, S_2, k)$: Improvements via Count-min sketch

1. Maintain two separate Count-min sketches for S_1, S_2
 2. Feed incoming k -mers to Count-min sketch while performing Reservoir Sampling in parallel. Update Count-min sketch accordingly.
 3. Once stream ends, Reservoir Sampling yields a k -mer σ and its position i . Let $S :=$ string in which σ is sampled, $S' :=$ the other string.
 4. Read the approximate frequency N of $N_k(\sigma, S)$ from Count-min sketch for S . Generate random integer m from range $[1, N]$.
 5. Read the approximate frequency N' of $N_k(\sigma, S')$ from Count-min sketch for S' . Return 1 or 0 based on whether $m \leq N'$.
-

Basically, rather than performing steps 2) and 3) by passing through the stream again, we shifted the burden to Count-min sketch. The “improvement” seems promising because if we assume that assume the Count-min sketch outputs the exact k -mer frequencies (i.e., $N = N_k(\sigma, S)$ and $N' = N_k(\sigma, S')$ in the algorithm). Then it turns out that the probability of f_{++} returning 1 is the same quantity as before:

$$\frac{2 \sum_{\tau} \min\{N_k(\tau, S_1), N_k(\tau, S_2)\}}{\sum_{\tau} N_k(\tau, S_1) + N_k(\tau, S_2)}$$

The current difficulty lies in measuring how the errors of Count-min sketch affects the true probability of f_{++} returning 1. This exploration is still ongoing and the author hopes to come up with explicit error bounds.