

Technical Report DIKU-TR-97/19
Department of Computer Science
University of Copenhagen
Universitetsparken 1
DK-2100 KBH Ø
DENMARK

September 1997

Floats, Integers, and Single Source Shortest Paths

Mikkel Thorup

Floats, Integers, and Single Source Shortest Paths

Mikkel Thorup

Department of Computer Science, University of Copenhagen
Universitetsparken 1, DK-2100 Copenhagen East, Denmark
mthorup@diku.dk, <http://www.diku.dk/~mthorup>

Abstract

Floats are ugly, but to everyone but theoretical computer scientists, they are the real thing. A linear time algorithm is presented for the undirected single source shortest paths problem with floating point weights.

1 Introduction

The technical goal of this paper is to present a linear time solution to the undirected single source shortest paths problem (USSSP) where the weights are floating points, or just floats. On a more philosophical level, the goal is to draw attention to the problem of making efficient algorithms for floats. Suppose, for example, we have an algorithm for the max-flow problem whose running time includes a factor $\log C$, where C is the maximal capacity. If we allow floating points, such an algorithm is not even polynomial, i.e. $\log C$ is the exponent of C , and the exponent is stored with $\log \log C$ bits.

Floating points are at least as used as integers, by everybody but theoretical computer scientists, who seem to prefer integers. To multiply two integers the faster way is often to convert them to double floats and send them to the floating point co-processor. Why? Because people simply don't care enough about integers to make an integer co-processor. Also, as theoretical computer scientists, we have to admit that floats do provide an elegant way of dealing with numbers in a large range. It is for good reasons that all other scientists and engineers have used them for centuries.

The rounding of floating point arithmetic is ugly in that, for example, addition is neither associative, nor commutative. Nevertheless, in this paper, we hope to indicate, that finding the structure of the rounding for a given problem, such as USSSP, can be an appealing challenge.

It should be noted that often floats do not cause any problems relative to integers. For problems like sorting, priority queues, and searching, there is no difference. The IEEE floating point standard is made such that interpreting the bit-string representation of floating points as representing integers, is order preserving. Hence we can feed floating points to an integer priority queue, if we just don't tell it that it is floats. Similarly, van Emde Boas' data structure [vBKZ77] works in time $O(\log \omega)$, where ω the the word length, no matter whether the words represent integers or floats.

Up to recently, all theoretical developments in the single source shortest paths problem (SSSP) were based in Dijkstra’s algorithm [Dij59], where vertices are visited in increasing order of distance to the source using a priority queue. Since the priority queue doesn’t care whether the input is integers or floats, all implementations of Dijkstra’s algorithm work equally well for integers and floats. However, recently, the author [Tho97] presented a linear time non-Dijkstra algorithm for undirected SSSP with integer weights. It is crucial to this algorithm that the weights are sorted with respect to their exponents. For integers in words, the exponent is at most ω , and then the exponents are easily sorted in linear time. However, for floating points, sorting the exponents is as hard as integer sorting, and this we do not know how to do in linear time.

Here, we show how to move edges around in the graph, preserving the floating point distances, ending with a series of independent subgraphs for which the exponents can be sorted in linear time. Given the sorted exponents, it is not too difficult to modify the algorithm from [Tho97] for integer USSSP to solve the floating point USSSP in linear time. Some remarks on this will be made in the journal version of [Tho97]. Given the sorted exponents, such a floating point version of the algorithm from [Tho97] may work very well in practice because of the efficiency of floating point arithmetic. However, here we give a self contained presentation, showing something in principle stronger; namely, that USSSP with floating point weights can be solved in linear time using a linear time oracle for USSSP with integer weights. Finally, we will make a few remarks on floats and max-flow.

2 Preliminaries

Our algorithm runs on a RAM, which models what we program in imperative programming languages such as C. The memory is divided into addressable words of length w . Addresses are themselves contained in words, so $\omega \geq \log n$. Moreover, we have a constant number of registers, each with capacity for one word. The basic assembler instructions are: conditional jumps, direct and indirect addressing for loading and storing words in registers, and some computational instructions, such as comparisons, addition, and multiplication, for numbers in registers. The space complexity is the maximal memory address used, and the time complexity is the number of instructions performed. All weights are floating point numbers, each contained in $O(1)$ words.

The numbers may be either integers, represented the usual way, or floating point numbers. A *floating point*, or just a *float*, is a pair $x = (e, m)$, where e is an integer, and m is a bit string $b_1 \cdots b_\varphi$. Then x represents the real number $2^e(1 + \sum_{i=1}^\varphi b_i/2^i)$. Thus $e = \lfloor \log_2 x \rfloor$. We call e the *exponent*, denoted $\mathbf{expo}(x)$, and m the *mantissa*, denoted $\mathbf{mant}(x)$. Often we will identify a float with the real number it represents. Both e and m are assumed to fit in a constant number of words. The number $\varphi = O(\omega)$ is fixed throughout a given computation and is referred to as the *precision*.

We let \oplus denote floating point addition. In this paper, for simplicity, when two floats x and y are added, they are always rounded *down* to nearest float (determined by the precision φ). Thus $x \oplus y \leq x + y$. This gives us the following basic rule:

$$\mathbf{expo}(x) > \mathbf{expo}(y) + \varphi \Rightarrow x \oplus y = y \oplus x = x \tag{1}$$

If instead we were rounding either up or down to nearest float, the rule would only apply if $\mathbf{expo}(x) > \mathbf{expo}(y) + \varphi + 1$, and all the calculations below, would have to be changed accordingly.

Let $G = (V, E)$, $|V| = n$, $|E| = m$, be an undirected connected graph with a distinguished source vertex s . Each edge $e \in E$ has a floating point weight $\ell(e)$ associated with it. The length of a path from s to some vertex v , is the floating point sum of the weights added up starting from s . More specifically, if the path is $P = (v_0, v_1, \dots, v_l)$, $s = v_0$, $v_l = v$, then the length of P is

$$((((\ell(v_0, v_1) \oplus (v_1, v_2)) \oplus \ell(v_2, v_3)) \cdots) \oplus \ell(v_{l-1}, v_l))$$

By $d(v)$ we denote the length of the shortest path from s to v . The floating point USSSP problem is that of finding $d(v)$ for all v . Different orders of adding the weights could give different answers, so it is natural to ask how big the errors can get. Let $d^*(v)$ denote the distance where the weights are summed using normal addition of the reals represented by the floats. Then $d^*(v)$ represents the ideal answer. Clearly $d^*(v) \leq d(v)$.

Observation 1 $\frac{d^*(v)-d(v)}{d(v)} \leq 2^{-\wp+\log_2 n}$

Proof: Set $e = \text{expo}(d(v))$. We are adding at most $n - 1$ numbers. Since the maximal exponent is $\leq e$, the maximal loss per number is $< 2^{e-\wp}$. Thus, the total error is $< (n - 1)2^{e-\wp} < 2^{e-\wp+\log_2 n}$. ■

In theory, since $\omega \geq \log_2 n$, we can always simulate $\log_2 n$ extra bit of precision without affecting the asymptotic running time. In practice, according to the IEEE standard format, with long floats, we have $\wp = 52$. Hence the relative error $\frac{d^*(v)-d(v)}{d(v)}$ is at most $2^{-52+\log_2 n}$, which is normally OK.

3 Sorting the exponents

We will now apply some different reductions to G constructing a graph G' so that the d -values of the nodes are unchanged (though some vertices may be identified), but so that for each component of $G' \setminus \{s\}$, the exponents vary by at most $n\wp$. Clearly the USSSP problem can be solved independently for each component of $G' \setminus \{s\}$, and the small variation in the exponents implies that they can be sorted in linear time.

Our first step is to construct a minimums spanning tree T for G in linear time [FW94]. Let $d_T(v)$ the weight of the path from s to v in T . Let $f_T(v)$ denote the weight of the first edge on this path. Clearly $d_T(v) \geq d(v)$. Also, since T is a minimum spanning tree, $f_T(v) \leq d(v)$.

Algorithm A: Reduce G by applying the following rules as long as possible.

1. If $v \neq s$, $\text{expo}(\ell(v, w)) > \text{expo}(d_T(v)) + \wp$ and $d_T(w) \geq d_T(v)$, replace (v, w) by an edge (s, w) of the same weight. If (v, w) was in T , it is replaced by (s, w) in T .

Since $d_T(v) \geq d(v)$, $d(v) \oplus \ell(v, w) = \ell(v, w)$. Also $d(w) \oplus \ell(v, w) \geq \ell(v, w) \geq d_T(v) \geq d(v)$. Thus no distances in G or T are changed.

2. If $(v, w) \in T$ and $\text{expo}(\ell(v, w)) < \text{expo}(f_T(v)) - \wp$, contract (v, w) .

Since $(v, w) \in T$, T remains a minimum spanning tree. Since $d(v) \geq f_T(v)$, $d(v) \oplus \ell(v, w) = d(v)$, so $d(w) \leq d(v)$. However, $f_T(v) = f_T(w)$, so symmetrically, $d(v) \leq d(w)$.

3. Remove any loop.

When no more rules apply, any edge (v, w) in resulting graph G' satisfies

(i) if $v \neq s$, $\text{expo}(\ell(v, w)) \leq \text{expo}(d_T(v)) + \wp$

(ii) $\text{expo}(\ell(v, w)) \geq \text{expo}(f_T(v)) - \wp$.

To see (ii), note that if it is not satisfied by a non-tree edge (v, w) , since T is a minimum spanning tree, (ii) should be satisfied by all the edges on the induced cycle, but these should have been contracted by rule 2, and then (v, w) should be a loop that should have been removed by rule 3.

We can now solve the USSSP problem for G , by solving it independently for each component of $G' \setminus \{s\}$. For simplicity, we assume that $G' \setminus \{s\}$ consists of exactly one component. Our first goal is to prove

Proposition 2 *If G' satisfies (i) and (ii) and $G' \setminus \{s\}$ is connected, the minimum and maximum exponent of an edge in G' differ by at most $n\wp$, where n is the number of nodes in G' .*

In order to prove Proposition 2, we first prove some technical lemmas. In each of the lemmas, we implicitly assume the conditions on G' from Proposition 2.

Lemma 3 *If v has depth $d > 0$, $\text{expo}(d_T(v)) \leq \text{expo}(f_T(v)) + (d - 1)\wp$.*

Proof: By induction on d . The statement is trivially true for $d = 1$. If $d > 1$ and u is the ancestor of v , by induction, $\text{expo}(d_T(u)) \leq \text{expo}(f_T(u)) + (d - 2)\wp = \text{expo}(f_T(v)) + (d - 2)\wp$. Applying (i), $\text{expo}(\ell(u, v)) \leq \text{expo}(d_T(u)) + \wp \leq \text{expo}(f_T(v)) + (d - 1)\wp$. Suppose $\text{expo}(\ell(u, v)) = \text{expo}(f_T(v)) + (d - 1)\wp$. Since $\text{expo}(d_T(u)) \leq \text{expo}(f_T(v)) - (d - 2)\wp$, $d_T(v) = d_T(u) \oplus \ell(u, v) = \ell(u, v)$. Since \oplus is increasing in both its arguments, it follows that $\text{expo}(d_T(v)) \leq \text{expo}(f_T(v)) + (d - 1)\wp$, as desired. ■

Lemma 4 *If $(v, w) \in G'$ where $v \neq s$, $\text{expo}(\ell(v, w)) \leq \text{expo}(f_T(v)) + \text{depth}_T(v)\wp$.*

Proof: By (i) and Lemma 3, $\text{expo}(\ell(v, w)) \leq \text{expo}(d_T(v)) + \wp \leq \text{expo}(f_T(v)) + \text{depth}_T(v)\wp$ ■

Lemma 5 *If (s, u) and (v, w) are in G' , $\text{expo}(\ell(v, w)) \leq \text{expo}(\ell(s, u)) + (n - 1)\wp$*

Proof: Let T_1, \dots, T_i be distinct subtrees of $T \setminus \{s\}$ such that $u \in T_1$, and for $i > 1$, there is an edge (x, y) from T_{i-1} to T_i . By induction on i , we will argue that if $(a, b) \in G'$ and $a \in T_i$,

$$\text{expo}(\ell(a, b)) \leq \text{expo}(\ell(s, u)) + \sum_{j=1}^i |V(T_j)|\wp. \quad (2)$$

Suppose $i = 1$. Since $\text{depth}_T(a) \leq |V(T_i)|$, (2) follows directly from Lemma 4.

Suppose $i > 1$. Let (x, y) be the edge entering T_i from T_{i-1} . By induction,

$$\text{expo}(\ell(x, y)) \leq \text{expo}(\ell(s, u)) + \sum_{j=1}^{i-1} |V(T_j)|\wp.$$

Since T is a minimum spanning tree $\ell(x, y) \geq f_T(y) = f_T(a)$. Thus, by Lemma 4,

$$\mathbf{expo}(\ell(a, b)) \leq \mathbf{expo}(f_T(a)) + \mathit{depth}_T(a)\varphi \leq \mathbf{expo}(x, y) + |V(T_i)|\varphi \leq \mathbf{expo}(\ell(s, u)) + \sum_{j=1}^i |V(T_j)|\varphi,$$

completing the proof of (2). The trees T_1, \dots, T_i can be chosen for any $(a, b) \in G'$, and since $\bigcup_j T_j \subseteq G' \setminus \{s\}$, the lemma follows. \blacksquare

Proof of Proposition 2: Let (s, u) be the minimum weight edge leaving s . Consider any edge (v, w) . By (2), $\mathbf{expo}(\ell(v, w)) \geq \mathbf{expo}(f_T(v)) - \varphi \geq \mathbf{expo}(\ell(s, u)) - \varphi$. At the same time, by Lemma 5, $\mathbf{expo}(\ell(v, w)) \leq \mathbf{expo}(\ell(s, u)) + (n - 1)\varphi$. Hence, exponents of weights in G' can vary by at most $n\varphi$. \blacksquare

Theorem 6 *We can sort the exponents of the edges in G' in linear time.*

Proof: First subtracting the minimum exponent from all other exponents, the maximum exponent becomes $n\varphi$, which is represented by $\log_2 n + \log_2 \varphi$ bits. If $\varphi \leq n$, each exponent is viewed as two $\leq \log_2 n$ bit characters, and we radix sort in two rounds, each round taking linear time.

In the extreme case where $\varphi > n$, since $\varphi = O(\omega)$, we have $\log_2 n + \log_2 \varphi = O(\omega/(\log n \log \log n))$. Then we can apply the linear time packed sorting from [AH92]. \blacksquare

4 Using an integer oracle

We will now solve the USSSP problem for the graph G' in Proposition 2 using the linear time integer USSSP oracle from [Tho97]. First we will estimate the exponents ± 1 . Second we will make exact calculations.

Distances with a bit too much precision

We will calculate the distances from s but sometimes using some extra precision. Since we always round down, extra precision means larger values. For each vertex v , we denote the obtained distance by $D(v)$. Referring to Observation 1, we get $d(v) \leq D(v) \leq d^*(v)$. Since $\varphi \geq \log_2 n$, we further get $D(v) \leq 2d(v)$. Consequently, $\mathbf{expo}(d(v)) \leq \mathbf{expo}(D(v)) \leq \mathbf{expo}(d(v)) + 1$.

Our first step is to subtract the exponent e of the smallest weight from all weight exponents (corresponding to division by 2^e). By Proposition 2, all weight exponents are now in the interval $[0, n\varphi)$.

We are going to proceed in rounds for $i = 1, \dots, n$. In round i , we are going to find the distances from s along paths where the maximal exponent of an edge weight is $< i\varphi$. We assume that this has already been done over paths where the maximal exponent is $< (i-1)\varphi$.

Consider the set $S = \{v \mid \mathbf{expo}(D(v)) < (i-1)\varphi\}$. Edges with exponents $\geq (i-1)\varphi$ are not going to give any better distances for vertices $v \in S$, so the distances to vertices in S may now be output. Next we *contract* S as follows. For all edges $(v, w) \in E$, $v \in S \setminus \{s\}$, $w \notin S$, we replace (v, w) by a *distance edge* (s, w) whose weight is $D(v) \oplus \ell(s, w)$. Afterwards $S \setminus \{s\}$ is removed. Clearly the remaining vertices have same distances in the reduced graph, and each edge (v, w) is only once replaced by a distance edge (s, w) .

Now contract all edges (v, w) with $\mathbf{expo}(\ell(v, w)) < (i-2)\wp$. This corresponds to reducing their weight to 0. We claim that this does not change the floating point distance to any vertex. By symmetry, it suffices to show that it does not change the distance to w . Suppose the shortest path to w goes through v . Since v was not contracted in S , $\mathbf{expo}(D(v)) \geq (i-1)\wp$, but this immediately implies $D(v) \oplus \ell(v, w) = D(w)$.

Here in round i , concerning the original edges, we restrict our attention to the surviving ones with exponent $< i\wp$. Since we have contracted all original edges with exponent $< (i-2)\wp$, this implies that edges will be considered in their original form for at most 3 rounds. Concerning distance edges, we restrict our attention to the surviving ones where the maximal edge weight on the corresponding path has exponent $< i\wp$. The weight of these distance edges is then $< (n-1)2^{i\wp} < 2^{(i+1)\wp}$. At the same time, we have contracted every distance edge whose exponent is $< (i-1)\wp$. Hence distance edges are considered for at most 3 rounds. In conclusion, each edge will be considered for a total of at most 6 times.

Before converting our floats to integers, we subtract $(i-3)\wp$ from every exponent, corresponding to dividing all weights by $2^{(i-3)\wp}$. Since the smallest weight exponent was $\geq (i-2)\wp$, it is now \wp , so all weights are now integers. All distances computed in round i , are based on original edges whose weight exponents are now $< i\wp - (i-3)\wp = 3\wp$, hence of weight $< 2^{3\wp}$. Thus, the maximal distance computed in round i is $< (n-1)2^{3\wp}$. Thus, we call our integer USSSP, where each integer is represented by $\lceil \log_2 n \rceil + 3\wp = O(\omega)$ bits. This takes time linear in the number of edges considered, and since each edge is considered in at most 6 rounds, the total running time over all rounds is linear in the total number of edges. The distances produced are essentially correct, except that in the calculations, we vary between \wp and $\log_2 n + 3\wp$ bits of precision, where we should really only have had a precision of \wp in every single addition. As pointed out above, this means that the exponents of the calculated distances are either correct, or at most one too large.

Exact results

In order to get the exact floating point distances $d(v)$, as defined in Section 2, we are going to proceed in rounds as above, but working with one exponent e at the time. We will go through the exponents e in increasing order. Recall that the exponents of the edge weights were sorted in Section 3. Inserting exponents of computed distances in the ordering is straightforward; for if $d(w) = d(v) \oplus \ell(v, w)$, $\mathbf{expo}(d(w)) \in \{\mathbf{expo}(d(v)), \mathbf{expo}(d(v)) + 1, \mathbf{expo}(\ell(v, w)), \mathbf{expo}(\ell(v, w)) + 1\}$. For each vertex v , we are going to start with our estimated distance $D(v) \geq d(v)$ computed in the previous subsection. We are going to decrease $D(v)$ to $d(v)$ over several rounds. Recall that when we start $\mathbf{expo}(d(v)) \leq \mathbf{expo}(D(v)) \leq \mathbf{expo}(d(v)) + 1$.

The goal of round e is to get $D(v) = d(v)$ for all v with $\mathbf{expo}(d(v)) \leq e$. Inductivity, we assume this has already been achieved for all v with $\mathbf{expo}(d(v)) \leq e-1$. Let S be the set of vertices v with $\mathbf{expo}(d(v)) \leq e-1$. These vertices are contracted, as in the previous section, that is, each outgoing edge (v, w) , $v \in S \setminus \{s\}, w \notin S$, is replaced by a distance edge (s, w) with $\ell(s, w) = d(v) \oplus \ell(v, w)$. Moreover, we set $D(w) = \min\{D(w), D(v) \oplus \ell(v, w)\}$.

Construct the graph G_e with the distance edges (s, v) with $\mathbf{expo}(D(v)) = e$, and all edges (v, w) with $\{\mathbf{expo}(D(v)), \mathbf{expo}(D(w))\} \subseteq \{e, e+1\}$. The latter includes all edges (v, w) with $\mathbf{expo}(d(v)) = \mathbf{expo}(d(w)) = e$. Subtract $e - \wp$ from all the exponents and convert to integers, rounding down, i.e. $x \mapsto \lfloor x/2^{e-\wp} \rfloor$. Call the integer USSSP algorithm, convert back to floats, and add $e - \wp$ to the exponents.

To see that the above is correct for any w with $\text{expo}(d(w)) = e$, consider a shortest path $v_0 \cdots v_l$ from $s = v_0$ to $w = v_l$. Let v_i be the last vertex with $\text{expo}(d(v_i)) < e$. By induction, $D(v_i) = d(v_i)$, and also, we contracted $v_i \in S$ correctly, setting $D(v_{i+1}) = D(v_i) \oplus \ell(v_i, v_{i+1}) = d(v_i) \oplus \ell(v_i, v_{i+1}) = d(v_{i+1})$. For $j = i+1, \dots, l-1$, we have $\text{expo}(d(v_j)) = \text{expo}(d(v_{j+1})) = e$, so (v_j, v_{j+1}) is an edge in G_e . Moreover, for each of the desired floating point additions $D(v_j) \oplus \ell(v_j, v_{j+1})$, the first term and the sum both have exponent e . This implies that the integer additions performed by the integer USSSP algorithm simulate the desired floating point additions with exactly the right precession φ . Thus, we end up with $D(w) = d(w)$, as desired.

Clearly, each edge is considered at most twice in its original form, and at most once as distance edge. Thus, the total time spent above is linear in the total number of edges.

Theorem 7 *There is a linear time Turing reduction from the floating point USSSP problem to the integer USSSP problem.*

Combining with the linear time algorithm for integer USSSP from [Tho97], we get

Corollary 8 *There is a linear time algorithm for USSSP problem with floating point weights.*

5 A remark on max-flow

Above it was shown how we can deal with floating points in connection with USSSP. In the introduction, we made some remarks concerning max-flow. We will now sketch how to replace an exponential factor $O(\log C)$ by a factor $O(\log n + \varphi)$. The conversion from integers to floats is a lot easier for max-flow than for single source shortest paths; for we are only interested in *one* flow value rather than a distance value for every single vertex.

Let f^* denote the (unknown) maximal flow value. First we find an s - t path with maximal minimal capacity D , that is, D is the maximal flow that can be pushed along a single path. Then $D \leq f^* \leq mD$. We can therefore reduce all capacities bigger than mD to mD without affecting the maximal flow. Concerning the small capacities, we can exploit that we are anyway going to accept an error in the order of $D/2^\varphi$. At most doubling this error, we can take all weights and round down to the nearest multiple of $D/(m2^\varphi)$. More precisely, we set $e = \text{expo}(D) - \lfloor \log_2 m \rfloor - \varphi$, subtract e from all exponents and round down to nearest integer. If C' is now the maximal capacity, $\log_2 C' \leq 2\lfloor \log_2 m \rfloor + \varphi = O(\log n + \varphi)$. Then we solve the max-flow problem with the reduced weights, and add e to the exponent of the resulting flow value.

References

- [AH92] S. Albers and T. Hagerup, Improved parallel integer sorting without concurrent writing, in *Proceedings of the 3rd ACM-SIAM Symposium on Discrete Algorithms*, pages 463–472, 1992.
- [Dij59] E.W. DIJKSTRA, A note on two problems in connection with graphs, *Numer. Math.* **1** (1959), 269–271.
- [FW94] M.L. FREDMAN AND D.E. WILLARD, Trans-dichotomous algorithms for minimum spanning trees and shortest paths, *J. Comp. Syst. Sc.* **48** (1994) 533–551.

- [Tho97] M. THORUP. Undirected Single Source Shortest Paths in Linear Time. To appear in *Proceedings of the 38th IEEE Symposium on Foundations of Computer Science, 1997*.
- [vBKZ77] P. VAN EMDE BOAS, R. KAAS, AND E. ZIJLSTRA, Design and implementation of an efficient priority queue, *Math. Syst. Th.* **10** (1977), 99–127.