
All-Pairs Shortest Paths in Spark

Charles Y. Zheng and Jinshu Wang

Department of Statistics
Stanford University
Stanford, CA 94305

{snarles, jinshuw}@stanford.edu

Arzav Jain

Department of Computer Science
Stanford University
Stanford, CA 94305

arzavj@stanford.edu

Abstract

We propose an algorithm for the All-Pairs-Shortest-Paths (APSP) problem suitable for implementation in Spark, and analyze its performance. We begin by considering distributed Floyd-Warshall, as proposed by Kumar and Singh (1991). Distributed Floyd-Warshall has asymptotically optimal scaling and can be implemented in Spark by using BlockMatrix to represent the APSP distance matrix. However, we observe that its implementation in Spark suffers from poor performance for medium-sized problems due the large number of global updates of the APSP distance matrix required for the algorithm. Since the lineage of the algorithm grows with the number of vertices n , it becomes necessary to use a proportional number of checkpoints which further impacts the efficiency of the algorithm. This motivates the consideration of an algorithm for APSP which requires fewer global update steps. We adapt an approach by Solomonik et al. (2013) based on the “divide and conquer” algorithm for APSP. Our algorithm reduces the number of global updates by a factor of b , where the block size b determines the amount of computation done in each iteration. By adjusting the block size b we obtain a favorable tradeoff between checkpointing costs and computation cost per iteration, resulting in far improved performance compared to Distributed Floyd-Warshall.

1 Summary

For the convenience of the grader we present an overview of our approach and our results. The rest of the paper gives a detailed explanation of the results in this section.

1.1 Problem Specification

Let $G = (V, E)$ be a graph with n vertices. Assume the input is given in the form of the adjacency matrix A of the graph stored as a BlockMatrix with equally sized square blocks. Specifically define the adjacency matrix A as a square matrix with dimension $n = |V|$, and entries

$$A_{ij} = \begin{cases} w_{i,j} & \text{if } (i \rightarrow j) \in E \\ 0 & \text{if } i = j \\ \infty & \text{if } (i \rightarrow j) \notin E \end{cases}$$

Let b be the size of the block, and let $n = b\ell$ so that ℓ^2 is the number of blocks. Write

$$A = \begin{pmatrix} A^{11} & A^{12} & \dots & A^{1\ell} \\ A^{21} & A^{22} & \dots & A^{2\ell} \\ \vdots & \vdots & \ddots & \vdots \\ A^{\ell 1} & A^{\ell 2} & \dots & A^{\ell \ell} \end{pmatrix}$$

so that A^{ij} is the (i, j) th block in the `BlockMatrix`.

The output is given by the APSP distance matrix S , where

$$S_{ij} = \begin{cases} \text{weight of shortest path} & \text{if there exists a path } i \rightarrow j \\ 0 & \text{if } i = j \\ \infty & \text{if there is no path } i \rightarrow j \end{cases}$$

Let S be stored as a `BlockMatrix` with the same dimensions and block sizes as A , so that

$$S = \begin{pmatrix} S^{11} & S^{12} & \dots & S^{1\ell} \\ S^{21} & S^{22} & \dots & S^{2\ell} \\ \vdots & \vdots & \ddots & \vdots \\ S^{\ell 1} & S^{\ell 2} & \dots & S^{\ell \ell} \end{pmatrix}$$

Let p be the number of workers. Let M be the memory of each worker. Suppose each worker holds K contiguous blocks. It must be the case that $K < M/b^2$. In fact, K is even smaller because each worker will have to hold additional data in memory.

1.2 Scaling

We consider the scaling $n \rightarrow \infty$ and $p \rightarrow \infty$. We do *not* assume a sparse graph G , so the number of edges can scale as $E \sim n^2$. However b must be constant since we assume each block must fit in memory.

1.3 Notation

Given an $n \times m$ matrix A and an $n \times m$ matrix B , define the entrywise minimum $C = \min(A, B)$ by

$$C_{ij} = \min(A_{ij}, B_{ij})$$

Meanwhile, given an $n \times k$ matrix A and a $k \times m$ matrix B , define the *min-plus* product $C = A \otimes B$ by

$$C_{i,j} = \min_{l=1}^k A_{il} + B_{lj}$$

for $i = 1, \dots, n$ and $j = 1, \dots, m$.

Define $\text{APSP}(A)$ as the all-pairs-shortest-distance matrix for adjacency matrix A . For example, $\text{APSP}(A)$ is obtained by running the Floyd-Warshall algorithm on A .

1.4 Algorithm

The algorithm consists of an *outer loop* with $\ell = n/b$ iterations. Each iteration culminates in the global update of the `BlockMatrix` S containing the intermediate values of the APSP distance matrix. Each outer loop iteration involves the execution of three distributed subroutines in sequence, called the A-step, the B-step and C-step. In addition, after every q iterations, the `BlockMatrix` S is checkpointed.

We first give a shorthand description of the algorithm without explicitly specifying the Spark operations used in each step or what data needs to be communicated at each step. In the analysis, we expand each step to describe the specific Spark operations needed, including the broadcasts, joins, etc. needed to transfer the necessary data across workers. Do note that $S^{(0)}, S^{(1)}, \dots, S^{(\ell)}$ refer to the sequence of `BlockMatrix` objects storing the results of each iteration.

Algorithm 1 Distributed Block APSP (shorthand)

```
function BLOCKAPSP(Adjacency matrix  $A$  given as a BlockMatrix with  $\ell$  row blocks and  $\ell$  column blocks)
   $S^{(0)} \leftarrow A$ 
  for  $k = 1, \dots, \ell$  do
    [A-step]
     $S^{kk(k)} \leftarrow \text{APSP}(S^{kk(k-1)})$ 
    [B-step]
    for  $i = 1, \dots, \ell, j = 1, \dots, \ell$  do in parallel
      if  $i = k$  and  $j \neq k$  then
         $S^{kj(k)} \leftarrow \min(S^{kj(k-1)}, S^{kk(k)} \otimes S^{kj(k-1)})$ 
      end if
      if  $i \neq k$  and  $j = k$  then
         $S^{ik(k)} \leftarrow \min(S^{ik(k-1)}, S^{ik(k)} \otimes S^{kk(k)})$ 
      end if
    end for
    [C-step]
    for  $i = 1, \dots, \ell, j = 1, \dots, \ell$  do in parallel
      if  $i \neq k$  and  $j \neq k$  then
         $S^{ij(k)} \leftarrow \min(S^{ij(k-1)}, S^{ik(k)} \otimes S^{kj(k)})$ 
      end if
    end for
    [D-step]
    if  $k \equiv 0 \pmod q$  then
      Checkpoint  $S^{(k)}$ 
    end if
  end for
  Return  $S = S^{(n/b)}$ , the APSP matrix in BlockMatrix form
end function
```

1.5 Optimality

The single-core cost of Floyd-Warshall, the best known single-core algorithm for APSP, is $O(n^3)$. A perfectly distributed form of Floyd-Warshall therefore has a total runtime of $O(n^3/p)$, in the asymptotic regime $n \rightarrow \infty$ and $p = O(n)$. Our algorithm achieves the same asymptotic runtime of

$$O\left(\frac{n^3}{p} + \frac{n}{b} + \right)$$

for details see section 3.

One can also consider the *communication cost* scaling in terms of the amount of data tranferred over the network. The paper by Solomonik et al. (2013) derived a theoretical lower bound on the communication cost of APSP as $\Omega(\frac{n^2}{p^{2/3}})$ words. In comparison, our algorithm communicates a total of $O(n^2\sqrt{p} + \frac{n}{b}\sqrt{p})$ words, which is worse by a power of p .

1.6 Communication Cost and Type

We analyze the *bandwidth* (total words sent) and the type of communication in each step of the algorithm. The A-step involves a one-to-one communication of a matrix of size $b \times b$ from a worker to the driver, involving a bandwidth of $O(b^2)$ words. The B-step involves a one-to-all broadcast of a matrix of size $b \times b$ from the driver to \sqrt{p} workers, hence a bandwidth of $O(b^2\sqrt{p})$ words. The C-step involves an all-to-all communication (a map-side join) where each worker recieves two $n/\sqrt{p} \times b$ matrices, and therefore entails a bandwidth of $O(nb\sqrt{p})$. Therefore the per-iteration bandwidth is $O((1 + \sqrt{p})b^2 + \frac{nb}{\sqrt{p}})$. The total bandwidth for the algorithm is $O(n^2\sqrt{p}b + n(1 + \sqrt{p})b)$. See section 3 for the derivation of the bandwidth per step.

2 Background

3 Analysis

In each step we give the *computational cost* (computation done on each worker), the *bandwidth* (number of words sent) and the total *runtime* of the entire step. Here we measure computational cost as runtime on a single machine, as distinguished from the distributed runtime as the time between the start of the algorithm and its termination as measured by an absolute clock. We give a non-asymptotic analysis given the following assumptions:

1. The time it takes to run Floyd-Warshall on a local matrix of size $b \times b$ is given by $\kappa_F b^3$
2. The time it takes to locally perform min-plus multiplication on matrices of size $a \times b$ and $b \times c$ be given by $\kappa_M ab^2 c$. The time to update $C \leftarrow \min(C, A \otimes B)$ is the same as the time to compute $A \otimes B$ since C can be modified in-place.
3. Separate the cost of communication and computation, so that sending messages and receiving messages is not included in the computational cost.
4. Time taken to send one message consisting of m words from one machine (whether a worker or driver) to another machine is

$$T(m) = \kappa_L + \kappa_T m$$

where κ_L is a latency constant describing the time it takes for data to travel through the network and κ_T is the transmission time per word.

5. Time taken to *broadcast* m words from the driver to p workers is given by

$$B(m, p) = \log(p)T(m)$$

due to the usage of bittorrent broadcast.

6. Each machine has a probability ϵ of instantaneously failing. The machine is instantly replaced, but all data in memory and disk is lost.
7. Let the time it takes to write to disk be given by $r\kappa_C m$, where κ_C is a constant for the disk write time and m is the amount of data each worker has to write. This is assuming the data is already on disk.
8. Assume no time cost for deleting data from disk.
9. We assume that the initial data A is stored on fault-tolerant backup, otherwise there is no guarantee that the job can complete. We assume a cost of $\kappa_B n^2$ to restore the cluster to the initial state from backup.

Note in particular that assumption 4 assumes a *uniform* network topology so the transmission rate from any worker to any other worker is equal. This is in contrast to the *grid* or *hypercube* topologies considered in most of the existing literature on distributed APSP.

Assumptions 6-9 deal with fault tolerance. In the following we will start by giving an analysis of *fault-free iterations*, and then as we analyze the global cost, we incorporate the extra cost from faulty iterations.

We do not account for random variation in the time taken to transmit messages. This is a potentially significant source of additional runtime in practical applications, but we are unaware of an accepted stochastic model of variation in transmission time, so we leave the modelling to future work. We also note that a rough approximation for the effect of variation in transmission time can be obtained by simply increasing the machine failure probability ϵ , since both effects have a similar dependence on the logarithm of the number of workers. Of course, it should be kept in mind that the analysis with the adjusted ϵ should only be used to compute runtime and not to optimize the frequency of checkpointing. Since checkpointing does not allieviate delays due to variation in transmission time, one should optimize checkpointing using the original ϵ representing machine failure rate.

3.1 A-step

In the k th iteration, the A-step is written in shorthand as $S^{kk(k)} \leftarrow \text{APSP}(S^{kk(k-1)})$. In fact, the updated `BlockMatrix` $S^{(k)}$ is not formed until the end of the C-step, and in the A-step $S^{kk(k)}$ only exists on the driver. $S^{ik(k)} \leftarrow \min(S^{ik(k-1)}, S^{ik(k)} \otimes S^{kk(k)})$ The detailed A-step is as follows:

1. The block $S^{kk(k-1)}$ is copied to the driver via invoking the `lookup` method with key (k, k) on the `BlockMatrix` $S^{(k-1)}$
2. The driver locally computes $S^{kk(k)} \leftarrow \text{APSP}(S^{kk(k-1)})$

Computation. Only the driver performs computation, running Floyd-Warshall on $S^{kk(k-1)}$. This costs $\kappa_F b^3$.

Bandwidth. The lookup involves a one-to-one communication between the worker holding S^{kk} and the driver. Exactly b^2 words are transmitted.

Runtime. The runtime consists of the time taken to transmit S^{kk} and the time to compute $\text{APSP} S^{kk}$. The transmission time is

$$\kappa_L + \kappa_T b^2$$

Therefore the total runtime is

$$\kappa_L + \kappa_T b^2 + \kappa_F b^3$$

3.2 B-step

In the k th iteration, the B-step is written in shorthand as $S^{kj(k)} \leftarrow \min(S^{kj(k-1)}, S^{kk(k)} \otimes S^{kj(k-1)})$ and

in parallel. This is accomplished by the following:

1. Form the `RDD rows` and the `RDD columns` by invoking the `filter` method on $S^{(k-1)}$ to keep only the blocks $S^{ik(k)}$ and $S^{kj(k)}$ for $i = 1, \dots, \ell$ and $j = 1, \dots, \ell$.
2. Broadcast $S^{kk(k)}$ from the driver to each worker holding a block in `rows` or `columns`. There should be \sqrt{p} workers storing blocks in `rows` and another \sqrt{p} workers storing blocks in `columns`.
3. Invoke `mapValues` on `rows` to perform the update $S^{kj(k)} \leftarrow S^{kk(k)} \otimes S^{kj(k-1)}$ and similarly on `cols` to perform the update $S^{ik(k)} \leftarrow S^{ik(k)} \otimes S^{kk(k)}$

Computation. Only the workers holding values in `rows` or `columns` perform computation, when `mapValues` is invoked. Each of those workers holds $\frac{n}{\sqrt{p}b}$ blocks. The update for each block entails a matrix multiply costing $\kappa_M b^3$, so the cost per worker is $\kappa_M \frac{nb^2}{\sqrt{p}}$, and the total cost is $2\kappa_M b^2 n$.

Bandwidth. The `broadcast` step entails sending a $b \times b$ matrix from the driver to $2\sqrt{p}$ workers. Therefore the bandwidth is $2\sqrt{p}b^2$

Runtime. The runtime cost of the broadcast is

$$\log(2\sqrt{p})(\kappa_L + \kappa_T b^2)$$

The runtime cost of computation is

$$\kappa_M \frac{b^2 n}{\sqrt{p}}$$

since each worker computes in parallel. Hence the total runtime is

$$\log(2\sqrt{p})(\kappa_L + \kappa_T b^2) + \kappa_M \frac{b^2 n}{\sqrt{p}}$$

3.3 C-step

In shorthand, the C-step is written as $S^{ij(k)} \leftarrow \min(S^{ik(k-1)}, S^{ik(k)} \otimes S^{kj(k)})$. Here it becomes important to note the partitioner used for the various RDDs. We assume that the same partitioner is used for all $S^{(0)}, \dots, S^{(\ell)}$ and for the RDDs `dupRows`, `dupCols`, and `temp` to be defined in the following.

1. Create RDD `dupRows` by invoking `flatMap` on `rows` to `flatMap` the key-value pair

$$(i, j, S^{ij})$$

to the key-value pairs

$$(1, j, S^{ij}), \dots, (\ell, j, S^{ij})$$

2. Create RDD `dupCols` in an analogous way
3. Create RDD `temp` by joining `dupCols`, `dupRows`, and $S^{(k-1)}$. The order of joining does not matter.
4. Create Blockmatrix $S^{(k)}$ by invoking `mapValues` to `temp`. The (i, j) entry of `temp` contains the blocks $S^{ij(k-1)}, S^{ik(k)}, S^{kj(k)}$ produced in the A-step and B-step. The effect of `mapValues` is to produce

$$S^{ij(k)} \leftarrow \min(S^{ik(k-1)}, S^{ik(k)} \otimes S^{kj(k)})$$

as (i, j) block in $S^{(k)}$.

Computation. The computation occurs when `mapValues` is called. There, each worker computes $S^{ij(k)} \leftarrow \min(S^{ik(k-1)}, S^{ik(k)} \otimes S^{kj(k)})$ for each of its $n^2/(b^2p)$ blocks. The computation per block is $\kappa_M b^3$, and hence the computation per worker is $\kappa_M \frac{n^2 b}{p}$. The total computational cost is $\kappa_M n^2 b$.

Bandwidth. Creation of `dupRows` and `dupCols` require bandwidth. Creation of `temp` and $S^{(k)}$ do not require bandwidth because of the map-side join. Each block in `rows` and `cols` must be sent to \sqrt{p} workers. Hence the bandwidth cost is $2b^2 \sqrt{p}$.

Runtime. Creation of `dupRows` requires each worker holding values in `rows` to transmit a message of size $\frac{nb}{\sqrt{p}}$ to \sqrt{p} other workers. This must be done sequentially under the current version of Spark (1.3.1) and hence takes time

$$\frac{p^2}{2} \kappa_L + \sqrt{p} \left(\kappa_T \frac{nb}{\sqrt{p}} \right) = p \kappa_L + \kappa_T nb$$

the term $\frac{p^2}{2} \kappa_L$ is the latency cost per worker of an all-to-all communication, since each pair of workers must communicate in order to organize the data transmission.

The same can be said for `dupCols`.

Creation of `temp` requires another pair of joins (but no other communication) and hence requires latency $\kappa_L p$.

In total, the runtime due to communication-related issues is

$$2\kappa_T nb + 2p\kappa_L$$

Adding the computation per worker, the total runtime is

$$2\kappa_T nb + 2p\kappa_L + \kappa_M \frac{n^2 b}{p}$$