

Web crawling: p. 72

provide one or more seed URLs as a string that will lead to many other web sites (tree) then links. obtain the content of these sites and store them.

How deep do you traverse: 2 or 3 deep depending on the application.

issues in web crawling: see p. 72

- An inverted index:
 - posting list one list for each term that appears in the collection.
- Posting list consists of individual postings:
 - each of which consists of document id, payload
- Payload may contain a variety of information
 - most commonly term frequency,
 - position of every occurrence of the term
 - other properties of the term: was it highlighted appeared in the title etc.
 - the list could be sorted by term frequency or some rank of the document.

term $\boxed{d_1 | p_1} \rightarrow \boxed{d_5 | p}$

term $\boxed{d_1 | p} \rightarrow \boxed{d_2 .}$

map (docid n, doc d)

$H \leftarrow$ new Association Array

for all terms $t \in \text{doc}$ do

$H\{t\} \leftarrow H\{t\} + 1$

for all terms $t \in H$ do

EMIT (term t , posting $\langle n, H\{t\} \rangle$)

reduce (term t , posting $\langle n, H\{t\} \rangle$ ^{list}
(n, f_1) (n, f_2) ...)

$P \leftarrow$ new List

for all posting p ^(a,f) in posting list

append (P , $\langle a, f \rangle$)

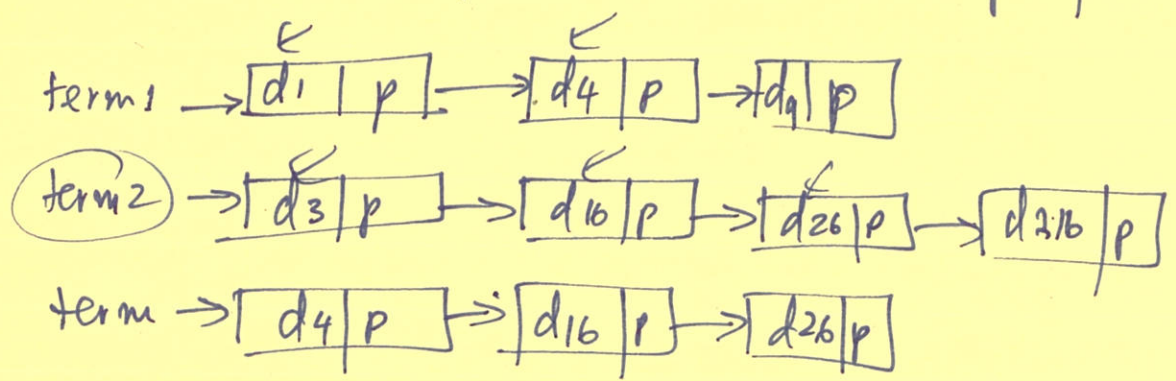
Sort P .

EMIT (term t , posting P)

Toy example

April 13, 2016
Spring 2016

CSE 487/587



p : payload

frequency # of times it occurs
 list of positions
 attributes of term

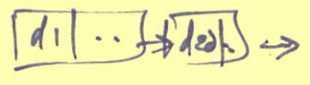
term → list is sorted by the document "rank"

```
class Mapper
  Mapper (docid n, doc d)
  method map (docid n, doc d)
  H ← new Associative Array
```

```
for all term t ∈ doc d do
  H {t} ← H {t} + 1
for all term t ∈ H do
  Emit (term t, posting ⟨n, H {t}⟩)
```

class Reducer

```
method reduce (term t, posting [⟨n, f1⟩ ... ⟨n, fn⟩])
  P ← new List
  for all posting ⟨n, f⟩ in posting do
    Append (P, ⟨n, f⟩)
  SORT (P)
  Emit (term t, posting P)
```



(2)

April 13, 2016

issue:

Sort is done by the reducer.

output is sorted by the "term" that is
the key

not by document.

solution:

key-value switch pattern

$\langle \langle \text{term } t, \text{doc } n \rangle, f(t) \rangle$

moving the "sort" out of the reducer.

(3)

April 13, 2014

doc1
one fish two fish

map

one → [d1 | 1]
fish → [d1 | 2]
two → [d1 | 1]

fish → [d1 | 2]
one → [d1 | 1]
two → [d1 | 1]

doc2
red fish blue fish

map

blue → [d2 | 1]
fish → [d2 | 2]
red → [d2 | 1]

doc3
one red bird

map

bird → [d3 | 1]
one → [d3 | 1]
red → [d3 | 1]

shuffle & Sort : aggregate ~~key by~~ values by key

reducer

bird → [d3 | 1]
blue → [d2 | 1]
fish → [d1 | 2] → [d2 | 2]
one → [d1 | 1] → [d3 | 1]
two → [d1 | 1] →
red → [d2 | 1] → [d3 | 1]

reduce is called how many times?
↳ distinct keys

April 13, 2016

Graph: $G(V, E)$

nodes edges

connections can be directed or undirected

social network

path planning

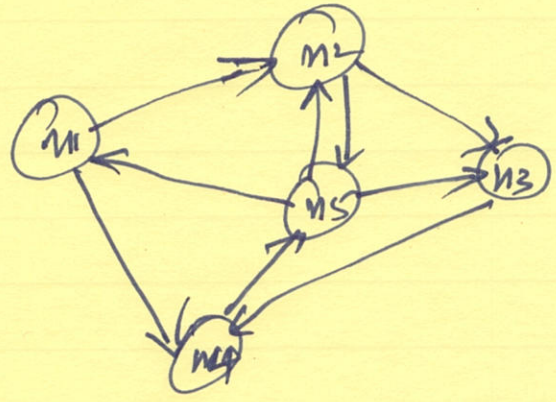
roadways

cable layout

Read the book

simple connected ~~graph~~ directed graph.

matrix: adjacency ^{matrix}



	n1	n2	n3	n4	n5
n1	0	1	0	1	0
n2	0	0	1	0	1
n3	0	0	0	1	0
n4	0	0	0	0	1
n5	1	1	1	0	0

representing

million nodes billion nodes

matrix : sparse matrix \Rightarrow lists

n1	[n2 n4]
n2	[n3 n5]
n3	[n4]
n4	[n5]
n5	[n1 n2 n3]

Breadth first search

adjacency lists

April 13, 2016

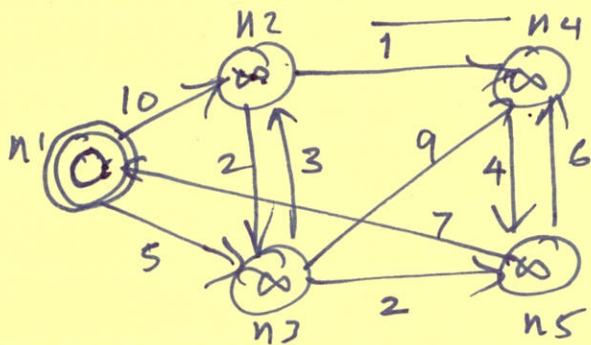
Dijkstra(G, w, s)

$d[s] = 0$; $d \rightarrow$ distance $Q[v, E]$
 for all vertex $v \in V$ (except first) $d[v] \leftarrow \infty$ are labels on the "nodes" computed by this algorithm
 $Q \leftarrow \{s\}$
 while ($Q \neq \text{empty}$)

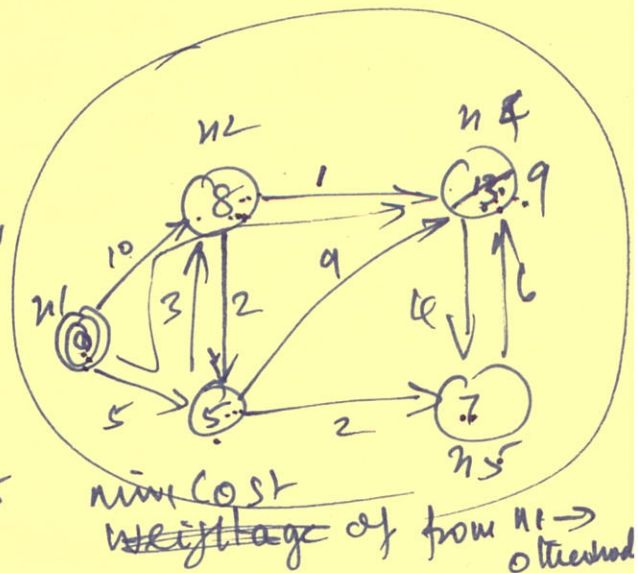
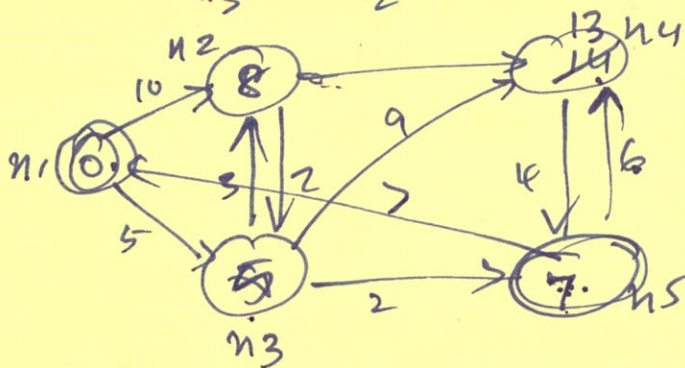
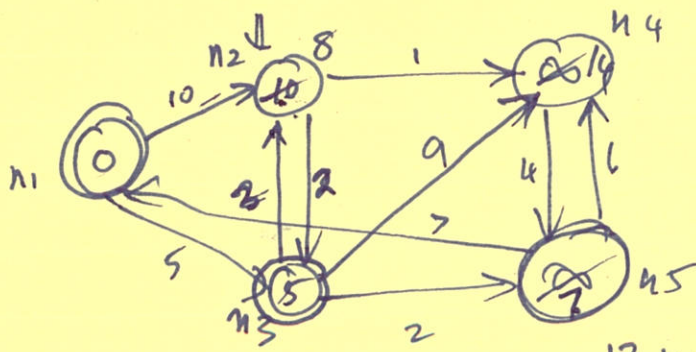
$u \leftarrow \text{ExtractMin}(Q)$;

for all vertex $v \in u$. adjacency list do
 if $d[v] > d[u] + w(u, v)$ then
 $d[v] \leftarrow d[u] + w(u, v)$

revisit



\Rightarrow n_1 starting point
 label all vertices with minimum weight from starting node



```

1: DIJKSTRA( $G, w, s$ )
2:    $d[s] \leftarrow 0$ 
3:   for all vertex  $v \in V$  do except source vertex
4:      $d[v] \leftarrow \infty$ 
5:    $Q \leftarrow \{V\}$ 
6:   while  $Q \neq \emptyset$  do
7:      $u \leftarrow \text{EXTRACTMIN}(Q)$ 
8:     for all vertex  $v \in u.\text{ADJACENCYLIST}$  do
9:       if  $d[v] > d[u] + w(u, v)$  then
10:         $d[v] \leftarrow d[u] + w(u, v)$ 

```

Figure 5.2: Pseudo-code for Dijkstra’s algorithm, which is based on maintaining a global priority queue of nodes with priorities equal to their distances from the source node. At each iteration, the algorithm expands the node with the shortest distance and updates distances to all reachable nodes.

As a refresher and also to serve as a point of comparison, Dijkstra’s algorithm is shown in Figure 5.2, adapted from Cormen, Leiserson, and Rivest’s classic algorithms textbook [41] (often simply known as *CLR*). The input to the algorithm is a directed, connected graph $G = (V, E)$ represented with adjacency lists, w containing edge distances such that $w(u, v) \geq 0$, and the source node s . The algorithm begins by first setting distances to all vertices $d[v], v \in V$ to ∞ , except for the source node, whose distance to itself is zero. The algorithm maintains Q , a global priority queue of vertices with priorities equal to their distance values d .

Dijkstra’s algorithm operates by iteratively selecting the node with the lowest current distance from the priority queue (initially, this is the source node). At each iteration, the algorithm “expands” that node by traversing the adjacency list of the selected node to see if any of those nodes can be reached with a path of a shorter distance. The algorithm terminates when the priority queue Q is empty, or equivalently, when all nodes have been considered. Note that the algorithm as presented in Figure 5.2 only computes the shortest distances. The actual paths can be recovered by storing “backpointers” for every node indicating a fragment of the shortest path.

A sample trace of the algorithm running on a simple graph is shown in Figure 5.3 (example also adapted from *CLR*). We start out in (a) with n_1 having a distance of zero (since it’s the source) and all other nodes having a distance of ∞ . In the first iteration (a), n_1 is selected as the node to expand (indicated by the thicker border). After the expansion, we see in (b) that n_2 and n_3 can be reached at a distance of 10 and 5, respectively. Also, we see in (b) that n_3 is the next node selected for expansion. Nodes we have already considered for expansion are shown in black. Expanding n_3 , we see in

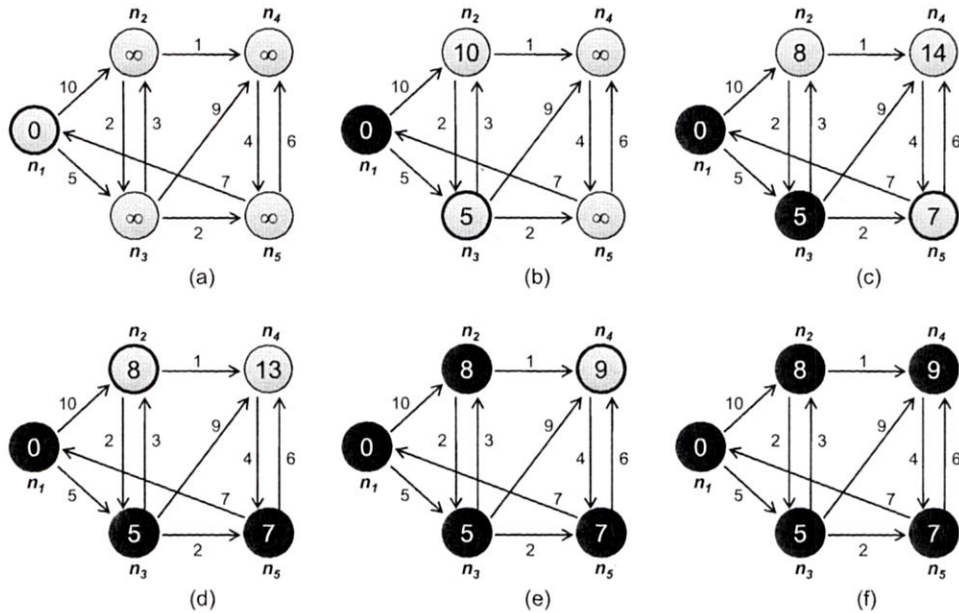


Figure 5.3: Example of Dijkstra's algorithm applied to a simple graph with five nodes, with n_1 as the source and edge distances as indicated. Parts (a)–(e) show the running of the algorithm at each iteration, with the current distance inside the node. Nodes with thicker borders are those being expanded; nodes that have already been expanded are shown in black.

(c) that the distance to n_2 has decreased because we've found a shorter path. The nodes that will be expanded next, in order, are n_5 , n_2 , and n_4 . The algorithm terminates with the end state shown in (f), where we've discovered the shortest distance to all nodes.

The key to Dijkstra's algorithm is the priority queue that maintains a globally-sorted list of nodes by current distance. This is not possible in MapReduce, as the programming model does not provide a mechanism for exchanging global data. Instead, we adopt a brute force approach known as parallel breadth-first search. First, as a simplification let us assume that all edges have unit distance (modeling, for example, hyperlinks on the web). This makes the algorithm easier to understand, but we'll relax this restriction later.

The intuition behind the algorithm is this: the distance of all nodes connected directly to the source node is one; the distance of all nodes directly connected to those is two; and so on. Imagine water rippling away from a rock dropped into a pond—that's a good image of how parallel breadth-first search works. However, what if there are multiple paths to the same node? Suppose we wish to compute the shortest distance

```

1: class MAPPER
2:   method MAP(nid  $n$ , node  $N$ )
3:      $d \leftarrow N.DISTANCE$ 
4:     EMIT(nid  $n$ ,  $N$ ) ▷ Pass along graph structure
5:     for all nodeid  $m \in N.ADJACENCYLIST$  do
6:       EMIT(nid  $m$ ,  $d + 1$ ) ▷ Emit distances to reachable nodes
1: class REDUCER
2:   method REDUCE(nid  $m$ , [ $d_1, d_2, \dots$ ])
3:      $d_{min} \leftarrow \infty$ 
4:      $M \leftarrow \emptyset$ 
5:     for all  $d \in \text{counts } [d_1, d_2, \dots]$  do
6:       if ISNODE( $d$ ) then
7:          $M \leftarrow d$  ▷ Recover graph structure
8:       else if  $d < d_{min}$  then ▷ Look for shorter distance
9:          $d_{min} \leftarrow d$ 
10:       $M.DISTANCE \leftarrow d_{min}$  ▷ Update shortest distance
11:      EMIT(nid  $m$ , node  $M$ )

```

Figure 5.4: Pseudo-code for parallel breadth-first search in MapReduce: the mappers emit distances to reachable nodes, while the reducers select the minimum of those distances for each destination node. Each iteration (one MapReduce job) of the algorithm expands the “search frontier” by one hop.

For more serious academic studies of “small world” phenomena in networks, we refer the reader to a number of publications [61, 62, 152, 2]. In practical terms, we iterate the algorithm until there are no more node distances that are ∞ . Since the graph is connected, all nodes are reachable, and since all edge distances are one, all discovered nodes are guaranteed to have the shortest distances (i.e., there is not a shorter path that goes through a node that hasn’t been discovered).

The actual checking of the termination condition must occur outside of MapReduce. Typically, execution of an iterative MapReduce algorithm requires a non-MapReduce “driver” program, which submits a MapReduce job to iterate the algorithm, checks to see if a termination condition has been met, and if not, repeats. Hadoop provides a lightweight API for constructs called “counters”, which, as the name suggests, can be used for counting events that occur during execution, e.g., number of corrupt records, number of times a certain condition is met, or anything that the programmer desires. Counters can be defined to count the number of nodes that have distances of ∞ : at the end of the job, the driver program can access the final counter value and check to see if another iteration is necessary.

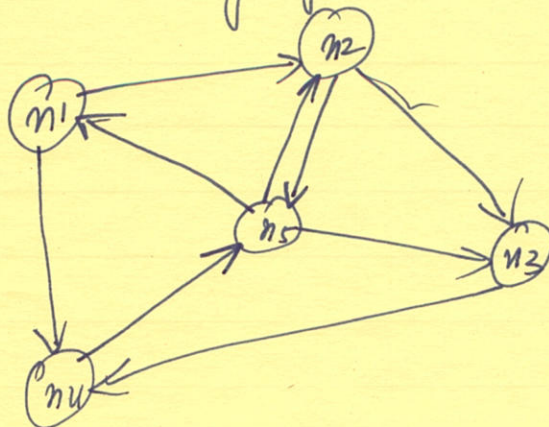
graphs are ubiquitous: ~~low~~ pathways, highways, friends, links of web sites.
flow of emails, messages

graphs are characterized by nodes and ~~set~~ edges.
Vertices links
connected: connection can be directed or undirected.

graphs are applied to many real-world problems:

- graph search and path planning
- graph clustering: identifying communities in social networks
- Minimum spanning trees: subset of the ~~of~~ edges that minimizes the sum of edge weights.
- Bipartite graph matching: matching two graphs: employers with potential employees
- Maximum flow:
 - identify special nodes: ~~single~~

simple directed graph:



	n_1	n_2	n_3	n_4	n_5
n_1	0	1	0	1	0
n_2	0	0	1	0	1
n_3	0	0	0	1	0
n_4	0	0	0	0	1
n_5	1	1	1	0	0

Dijkstra's algorithm:

Dijkstra(G, w, s) ^{source node}

$d(s) \leftarrow 0$

for all vertex $v \in V$ do

$d(v) \leftarrow \infty$

$Q \leftarrow \{v\}$

~~while~~ while $Q \neq \emptyset$ do

$u \leftarrow \text{ExtractMin}(Q)$

for all vertex $v \in u$. AdjacentList do

if $d[v] > d[u] + w(u, v)$ then

$d[v] \leftarrow d[u] + w(u, v)$

Key:

Priority queue of nodes sorted by distances

$G(V, E)$
 s source node
 w
 $w(u, v) \geq 0$

Base case algorithm: +1 source node \rightarrow adjacent node

