

# Week 09 Graphs

## SSSP for weighted and unweighted graphs

MIDS 261

Notebook:

[https://www.dropbox.com/s/zjm0ka9jr7yvlih/demo9\\_workbook\\_MASTER\\_2023\\_07\\_04.ipynb?dl=0](https://www.dropbox.com/s/zjm0ka9jr7yvlih/demo9_workbook_MASTER_2023_07_04.ipynb?dl=0)

# Housekeeping

---



Home  
Announcements  
Syllabus  
Modules  
Grades  
Files

Courses

Calendar  
Inbox  
History  
Placement Portal  
Help

- Once teams are formed, everyone needs to fill in their team details on DC under the "People" section. You can join a group by adding your name to one of the final project groups in DC. Click the "People" TAB on the left-hand navigation menu, and select the "FP\_SectionK\_GroupN" tab to view available groups. Each group is limited to a max of four students inside DC.
- Please see [here](#) to learn more about joining a group in DC (do not create your own student groups, you must join a pre-setup group created by your Instructors for this project).
- Please use the team following team name prefix, FP\_SectionK\_GroupN\_ as is but feel free to add a team suffix such as RootFinders to yield a team name of FP\_Section1\_Group2-RootFinders; when K denotes the section you belong to and N denotes a unique group number with your section.

For more details on the final project, please see [Module - Final Project - Flight Delays](#).

## Schedule at a glance

Phase Description (For details please see detailed phase descriptions and deliverables)	Submission details (or just see <a href="#">Module - Final Project - Flight Delays</a> )	Week number	DataBricks cluster details
Phase 0: Finalize Teams (teams of four people)	Due Sunday of Week 9	Week 9	Available from Monday/Tues week 10 (One master and one worker)
Phase I - project plan, describe datasets, joins, tasks, and metrics	Phase 1 Update: Notebook. Due Sunday of Week 10	Week 10	
Phase II - EDA, Baseline Pipeline on all data: Scalability, Efficiency, Distributed/parallel Training, and Scoring Pipeline	Phase 2 Update: Video (2 minutes max) + Notebook. Due Sunday of Week 11	Week 11	AutoScaling enabled from Monday of this week (up to 5)
Phase III Feature engineering + hyperparameter tuning, + in-class review	Phase 3 update: Notebook. Due by Sunday of Week 12	Week 12	
Phase IV - Advanced model architectures and loss functions, select an optimal algorithm, fine-tune & Final report write up	Phase 4 update: 2-minute Video + Notebook due by Sunday of Week 13	Week 13	AutoScaling enabled from Monday of this week (up to 10 depending on budget)  NOTE: Clusters will be scaled back on Sunday of this week.
Phase V - In-class Presentation	Slides uploaded prior to in-class presentation.	Week 14	



Home

Announcements

Syllabus

Modules

Grades

Files

Zoom Live Sessions

People

Collaborations

Faculty Course  
Community

Pages

BigBlueButton

Quizzes

Outcomes

Discussions

Assignments

Rubrics

Settings

Everyone

Student Groups

Project Groups

+ Group Set

Self sign-up is enabled for these groups. [?](#)

Groups are limited to 4 members.

+ Import

+ Group

⋮

### Unassigned Students (63)

Search users

⋮ Anshuman Awasthi +

⋮ Ahmad Azizi +

⋮ Rohit Bakshi +

⋮ Job Eliezek Bangayan +

⋮ Noriel Bargo (*He/Him*) +

⋮ Andrew Beckerman +

⋮ Nicholas Brathwaite +

⋮ Aris Chalini +

⋮ Evan Chan +

⋮ Justin Chan +

⋮ Sitao Chen +

⋮ Tyler Chi +

⋮ Nathan Chiu +

### Groups (28)

▶ FP\_Section1\_Group1

0 / 4 students

⋮

▶ FP\_Section1\_Group2

Full  
4 / 4 students

⋮

▶ FP\_Section1\_Group3

Full  
4 / 4 students

⋮

▶ [FP\\_Section1\\_Group4](#)

Full  
4 / 4 students

⋮

▶ FP\_Section2\_Group1

2 / 4 students

⋮

▶ FP\_Section2\_Group2

Full  
4 / 4 students

⋮



Account



Dashboard



Courses



Calendar



Inbox



History

Placement  
Portal

Help

Home

Everyone

Student Groups

Project Groups

+ Group Set

Announcements

+ Import

+ Group

⋮

Syllabus

Modules

Grades

Files

Zoom Live Sessions

People

Collaborations

Faculty Course  
Community

Pages

BigBlueButton

Quizzes

Outcomes

Discussions

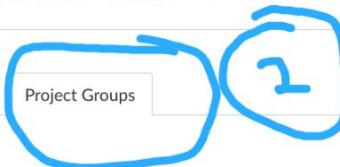
Assignments

Rubrics

Google Drive

Ally Course  
Accessibility Report

Settings



Unassigned Students (56)

Groups (8)

1. Click on 'Project Groups'

2. Click on 'FP\_Section1\_Group 1'

3. Pick a group

Drag your name to your group

Group	Students	⋮
FP_Section1_Group 1	0 / 4 students	⋮
FP_Section1_Group 2	0 / 4 students	⋮
FP_Section1_Group 3	0 / 4 students	⋮
FP_Section1_Group 4	0 / 4 students	⋮
FP_Section2_Group 1	0 / 4 students	⋮
FP_Section2_Group 2	0 / 4 students	⋮
FP_Section2_Group 3	0 / 4 students	⋮
FP_Section2_Group 4	0 / 4 students	⋮



Account



Dashboard



Courses



Groups



Calendar



Inbox



History



Placement  
Portal



Help

Home

Everyone

Groups

+ Group



Announcements

Syllabus

Search Groups or People

Modules

Grades

Files

Zoom Live Sessions

People

Collaborations

Google Drive

FP\_Section1\_Group 1 FP\_SectionK\_Group [Visit](#)

2 students

[Leave](#)

FP\_Section1\_Group 2 FP\_SectionK\_Group

0 students

[Switch To](#)

FP\_Section1\_Group 3 FP\_SectionK\_Group

0 students

[Switch To](#)

FP\_Section1\_Group 4 FP\_SectionK\_Group

0 students

[Switch To](#)

FP\_Section2\_Group 1 FP\_SectionK\_Group

0 students

[Switch To](#)

FP\_Section2\_Group 2 FP\_SectionK\_Group

0 students

[Switch To](#)

---

# Quick review

9: Graph Algorithms at Scale	
	Week 09 Async Slides Combined.pdf
	Weekly Introduction 9
	Networks Introduction and Motivation
	Applications Graphs
	Graph Definitions
	Shortest Path Introduction
	Single-Source Shortest-Path Unweighted Graphs
	BFS for Unweighted Graphs: Directed and Animations
	BFS for Unweighted Graphs: Undirected and Animations
	SSSP for Weighted Graphs BFS
	Quiz 9-1
	Quiz
	Question: Quiz 9-1
	Quiz
	SSSP for Weighted Graphs: Dijkstra's Algorithm
	Distributed SSSP: Unweighted Graph
	Distributed SSSP: Unweighted Graph Algorithm
	Distributed SSSP: Weighted Graph
	Question: Quiz 9-2
	HW5 - PageRank Jul 23   100 pts

# Some common problems

Traverse all the vertices in a graph (Keys and Rooms in LeetCode)

Whether two vertices are connected or not (Division Evaluation)

Shortest path between two vertices (Shortest Path in Binary Matrix, Word Ladder I & II, Walls and Gates)

Number of connected components (Number of Islands)

Dependency path (Course Schedule II, Alien Dictionary)

Vertex grouping (Is Graph Bipartite?)

Detect cycle in a graph (Course Schedule I, Valid Tree Graph)

All possible paths (Word Search, Binary Tree Paths)

## 841. Keys and Rooms

Medium

550

44

Favorite

Share

There are  $N$  rooms and you start in room  $0$ . Each room has a distinct number in  $0, 1, 2, \dots, N-1$ , and each room may have some keys to access the next room.

Formally, each room  $i$  has a list of keys  $\text{rooms}[i]$ , and each key  $\text{rooms}[i][j]$  is an integer in  $[0, 1, \dots, N-1]$  where  $N = \text{rooms.length}$ . A key  $\text{rooms}[i][j] = v$  opens the room with number  $v$ .

Initially, all the rooms start locked (except for room  $0$ ).

You can walk back and forth between rooms freely.

Return `true` if and only if you can enter every room.

### Example 1:

**Room 0 1 2 3**

**Input:** `[[1],[2],[3],[]]`

**Output:** `true`

**Explanation:**

We start in room  $0$ , and pick up key  $1$ .

We then go to room  $1$ , and pick up key  $2$ .

We then go to room  $2$ , and pick up key  $3$ .

We then go to room  $3$ . Since we were able to go to every room, we return `true`.

Leetcode

### Example 2:

**Input:** `[[1,3],[3,0,1],[2],[0]]`

**Output:** `false`

**Explanation:** We can't enter the room with number  $2$ .

## 127. Word Ladder

Medium

1749

882

Favorite

Share

Given two words (*beginWord* and *endWord*), and a dictionary's word list, find the length of shortest transformation sequence from *beginWord* to *endWord*, such that:

1. Only one letter can be changed at a time.
2. Each transformed word must exist in the word list. Note that *beginWord* is *not* a transformed word.

**Note:**

- Return 0 if there is no such transformation sequence.
- All words have the same length.
- All words contain only lowercase alphabetic characters.
- You may assume no duplicates in the word list.
- You may assume *beginWord* and *endWord* are non-empty and are not the same.

**Example 1:**

**Input:**

```
beginWord = "hit",
endWord = "cog",
wordList = ["hot","dot","dog","lot","log","cog"]
```

**Output:** 5

**Explanation:** As one shortest transformation is "hit" -> "hot" -> "dot" -> "dog" -> "cog", return its length 5.

## 269. Alien Dictionary

There is a new alien language which uses the latin alphabet. However, the order among letters are unknown to you. You receive a list of words from the dictionary, where words are sorted lexicographically by the rules of this new language. Derive the order of letters in this language.

For example, Given the following words in dictionary,

```
[  
    "wrt",  
    "wrf",  
    "er",  
    "ett",  
    "rftt"  
]
```

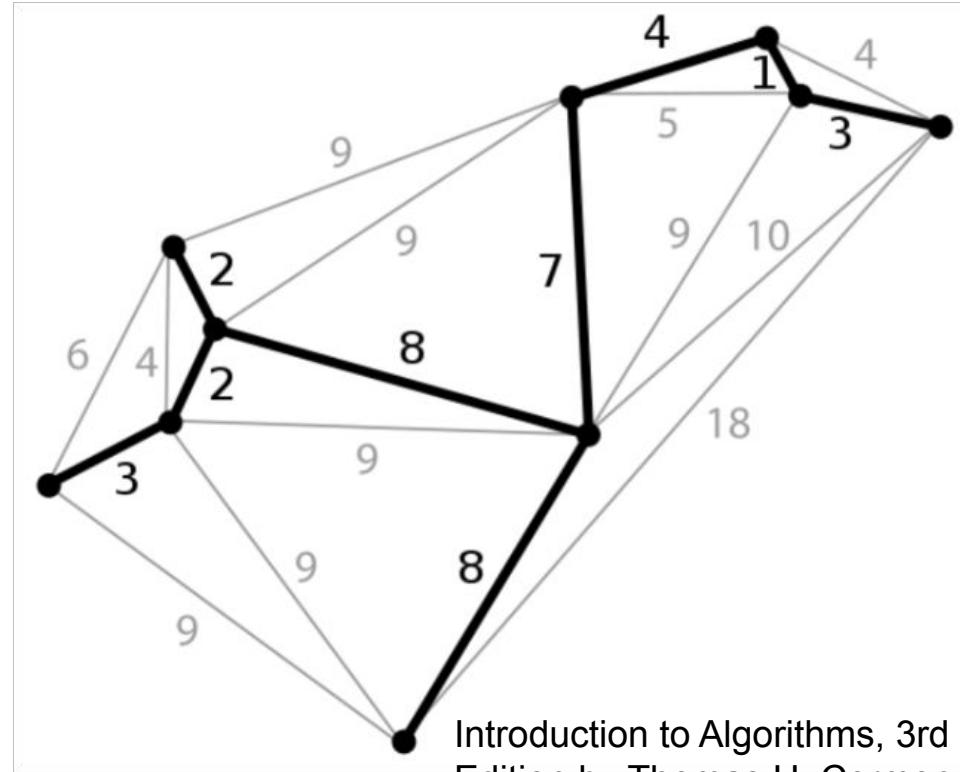
The correct order is: "wertf".

Note: You may assume all letters are in lowercase. If the order is invalid, return an empty string. There may be multiple valid order of letters, return any one of them is fine.

# Minimum Spanning Tree

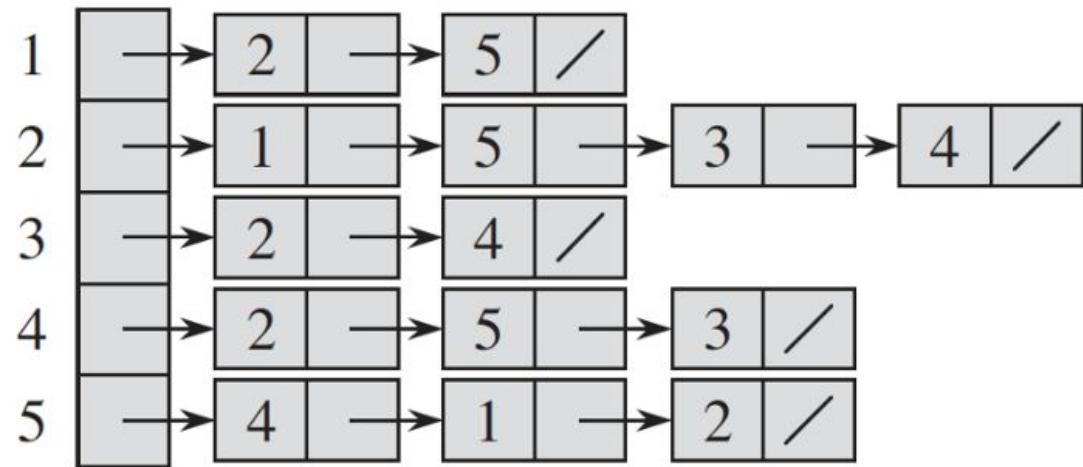
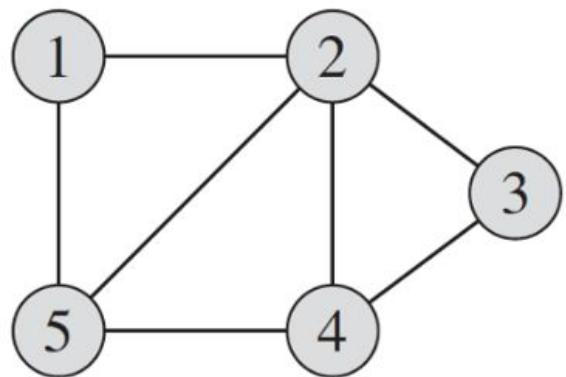
MST is a subset of the edges that connects all the vertices of a connected, edge-weighted undirected graph, without any cycles and with the minimum possible total edge weight.

Telco/Cable/Pipeline problem



Introduction to Algorithms, 3rd  
Edition by Thomas H. Cormen,  
Charles E. Leiserson, Ronald L.  
Rivest, Clifford Stein

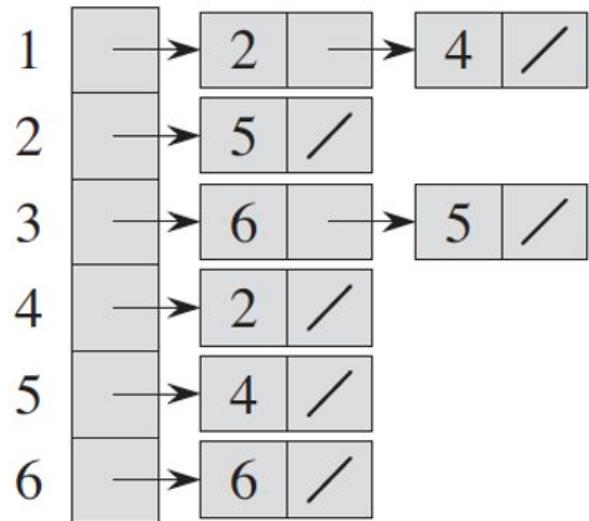
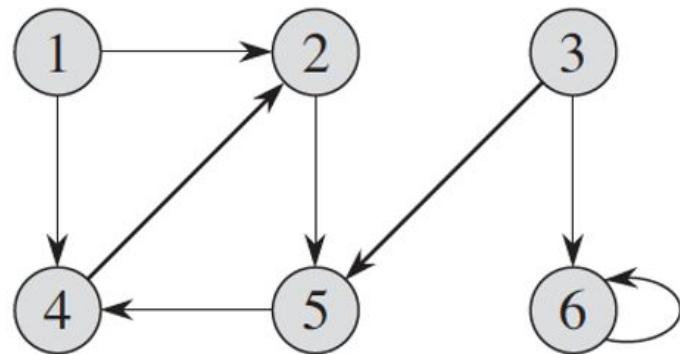
# Adjacency List (Undirected Graph)



Introduction to Algorithms, 3rd Edition by Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, Clifford Stein

# Adjacency List (directed Graph)

---



---

Introduction to Algorithms, 3rd Edition by Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, Clifford Stein

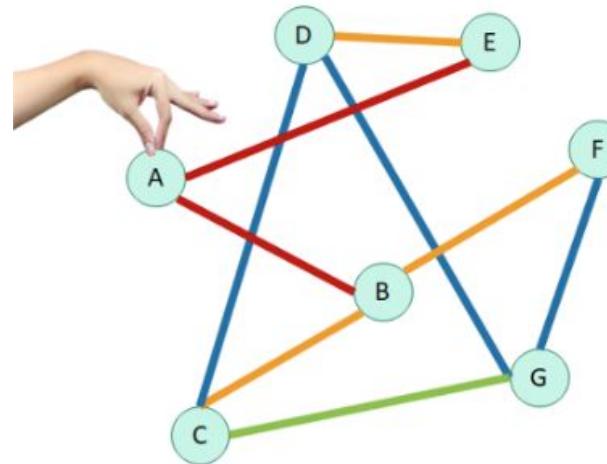
---

# BFS Traversal of a graph

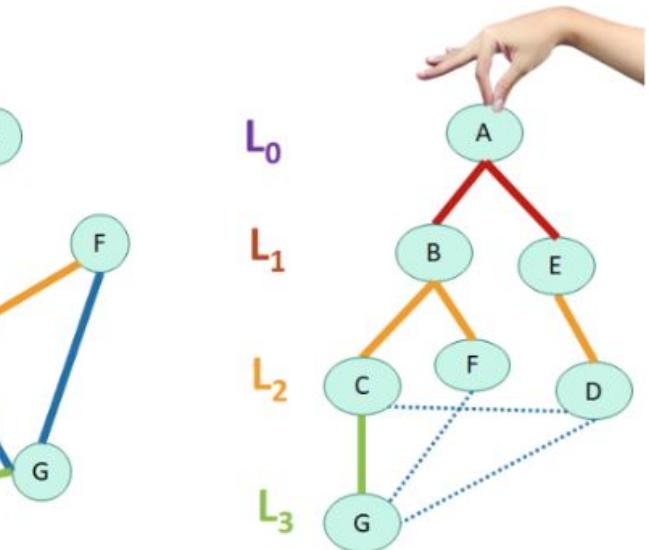
# BFS shortest path for unweighted graphs

BFS traverses by level

Number of levels is  
also the shortest  
distance between two  
vertices



Original Graph

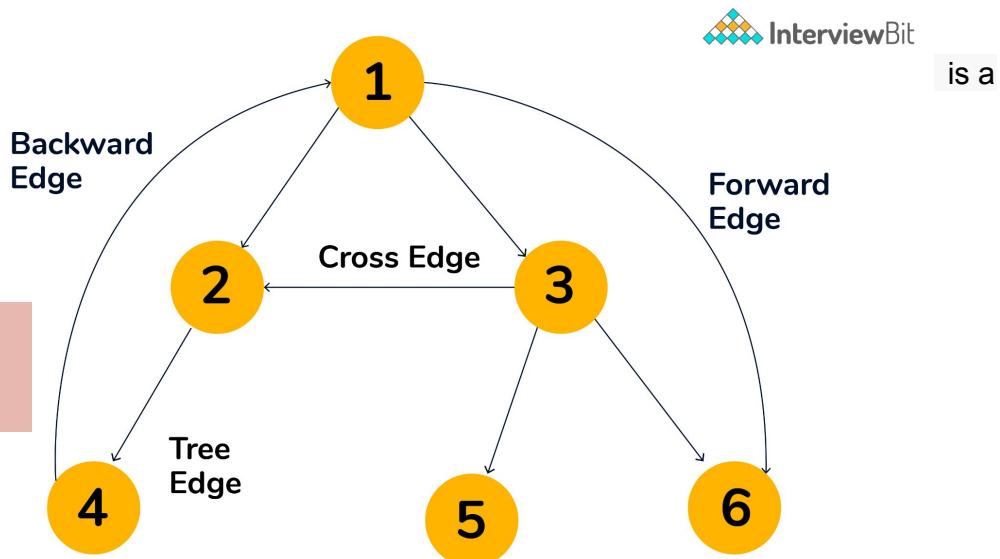


Visualized  
BFS Traversal

# Cross edges, back edges (loops in graphs)

- Tree Edge: The edge which is present in the tree obtained after applying DFS on the graph. All the Green edges are tree edges.
- Forward Edge: Edge  $(u, v)$  such that  $v$  is descendant but not part of the DFS tree. Edge from node 1 to node 6 is a forward edge.
- Back edge:
  - Edge  $(u, v)$  such that  $v$  is ancestor of  $u$ . Edge from node 3 to node 2 is a back edge.
  - Presence of back edge indicates a cycle in the graph.
- Cross Edge: It is an edge which connects two nodes that have no parent-child relationship between them. Edge from node 3 to node 4 is a cross edge.

Visited arrays prevent loops in BFS traversal and SSSP for unweighted graphs



<https://www.interviewbit.com/courses/programming/graph-data-structures-algorithms/breadth-first-search/>

1. MASTER Unit 9 - Graph Algorithms at Scale  
1.0.1. Notebook Set-Up

2. Exercise 1. Graphs Overview  
3. Warm-up questions:  
4. Exercise 2. Data Structures Review.

5. Exercise 3. Graph Traversal.

5.1. Depth-First Search

5.2. Breadth-First Search

5.3. Other great algorithm books:

5.3.1. BFS traversal

5.3.2. DFS

5.4. BFS versus DFS (no animation)

5.5. Example random graph

6. TASK: Write BFS Search Function In Apache Spark And Pandas Dataframe:

6.1. TODO: iterative search over undirected graph

7. Exercise 4. SSSP for unweighted graphs.

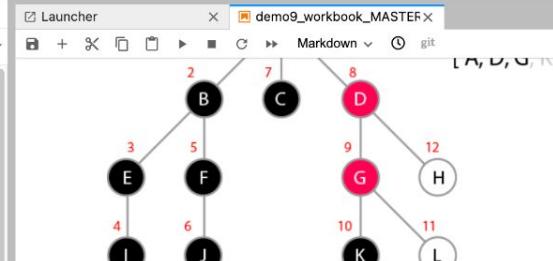
7.1. SSSP with a more illustrative example

7.2. SSSP for a random graph

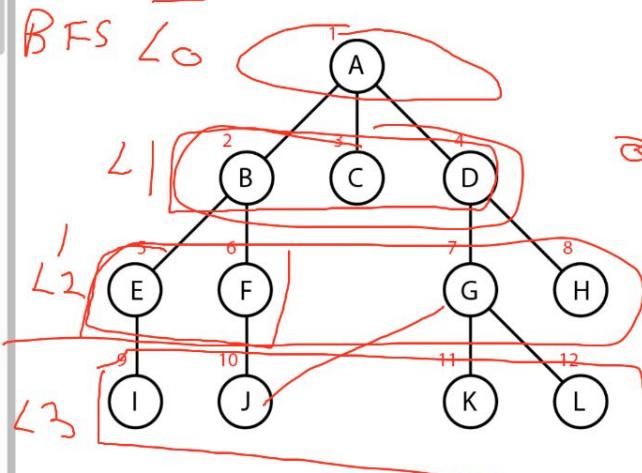
8. Exercise 5: Dijkstra's algorithm (SSSP for weighted graphs)

8.0.0.1. What about longest paths?

8.1. 8.1. ¶



#### 5.4. BFS versus DFS (no animation)



STATE of nodes

Visited = [A B C. ]

Queued [ ... - ]

Unvisited [ -- - ]

Breadth First - FIFO

Start with node A and put it on the Queue.

Pop the "first node", mark it as "visited", and put all unvisited reachable nodes on the queue:

state of f	Q node	neighbors	State	Dist
[A]	A	B C D	Q	0
[B, C, D]		E F	V	1
[C, D, E, F]		G H	V	2
[D, E, F]		I J	V	3
[E, F, G, H]		K L	V	4
[F, G, H, I]				
[G, H, I, J]				
[H, I, J, K, L]				
[I, J, K, L]				
[J, K, L]				
[K, L]				
[L]				
[]				

Single shortest path calculation

Distance and levels: when we put each neighbor in the Queue we add (neighbor, distance to source node) structure to the Queue instead

```

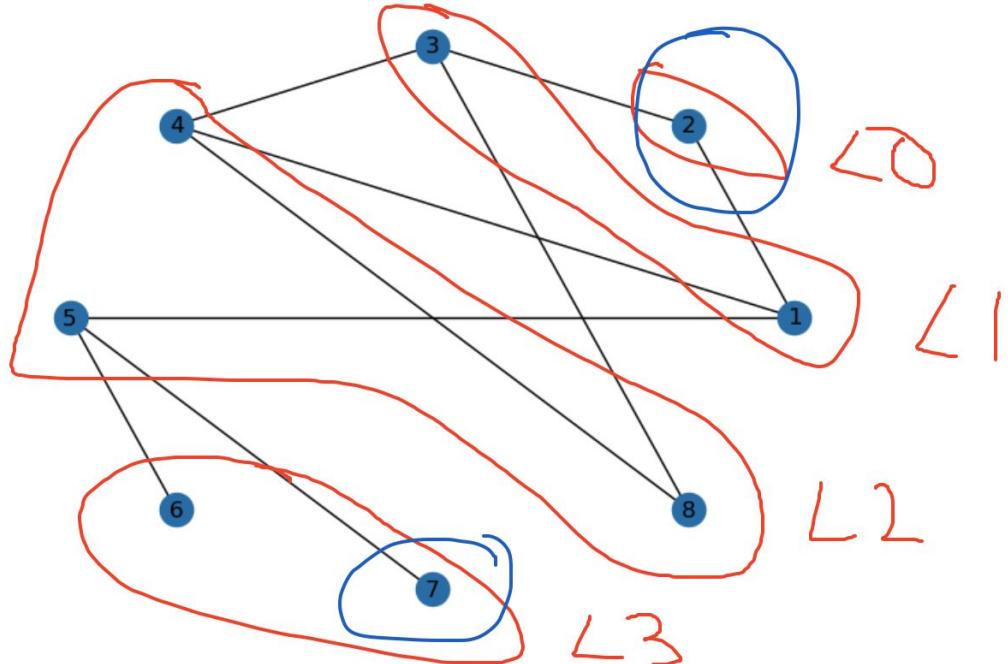
4.1 print(bfs_shortest_path(ADJ_GRAPH, 2, 7) is bfs_shortest_path(ADJ_GRAPH_2, 2, 7))
[1, 2, 3, 4, 5, 6, 7, 8]
[(1, 2), (1, 4), (1, 5), (2, 3), (3, 4), (3, 8), (4, 8), (5, 6), (5, 7)]
[2, 4, 5]
3
{1: [2, 4, 5], 2: [1, 3], 3: [2, 4, 8], 4: [3, 1, 8], 5: [1, 6, 7], 6: [5], 7: [5], 8: [4, 3]}
BFS Shortest paths

```

```

visited: []
queue: deque([[2]])
visited: [2]
queue: deque([[2, 1], [2, 3]])
visited: [2, 1]
queue: deque([[2, 3], [2, 1, 2], [2, 1, 4], [2, 1, 5]])
visited: [2, 1, 3]
queue: deque([[2, 1, 2], [2, 1, 4], [2, 1, 5], [2, 3, 2], [2, 3, 4], [2, 3, 8]])
visited: [2, 1, 3]
queue: deque([[2, 1, 4], [2, 1, 5], [2, 3, 2], [2, 3, 4], [2, 3, 8]])
visited: [2, 1, 3, 4]
queue: deque([[2, 1, 5], [2, 3, 2], [2, 3, 4], [2, 3, 8], [2, 1, 4, 3], [2, 1, 4, 1], [2, 1, 4, 8]])
bfs_shortest_path(ADJ_GRAPH, 2, 7) is ([2, 1, 5, 7], 3)

```

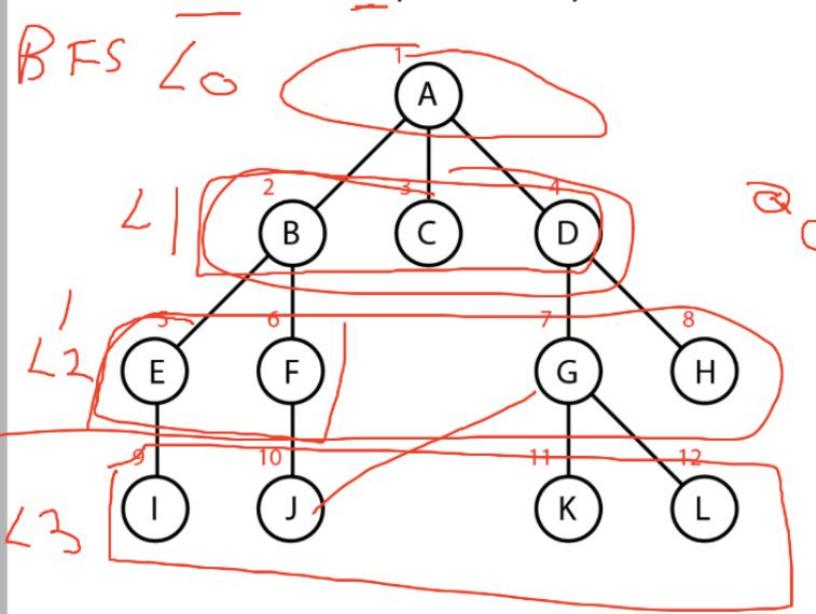


---

# **SSSP of an unweighted graph**

# BFS traversal with distances (via BFS tree)

- Each node is in a single state:
  - Queue
  - Visited (to avoid loops)
  - Unvisited (implicite)



```
[8]: 1 from collections import deque
2 # Adapted from: [here](https://www.geeksforgeeks.org/shortest-path-unweighted-graph/)
3 def bfs_shortest_path(graph, start):
4     # Create a queue and add the starting vertex to it
5     queue = deque([start])
6
7     # Create an array to keep track of the distances from the starting vertex to all other vertices
8     distances = [float('inf')] * len(graph)
9     distances[start] = 0
10
11    # Create a set to keep track of visited vertices
12    visited = set()
13
14    # Perform BFS
15    while queue:
16        # Dequeue the next vertex
17        vertex = queue.popleft()
18        visited.add(vertex)
19
20        # Update the distances of neighbors
21        for neighbor in graph[vertex]:
22            if neighbor not in visited:
23                distances[neighbor] = distances[vertex] + 1
24                queue.append(neighbor)
25
26    return distances
27
28
29 # Example graph: unweighted, directed graph with 5 vertices
30 # Vertices are represented by integers 0 through 4
31 # Edges: (0, 1), (0, 2), (1, 2), (1, 3), (2, 3), (3, 4)
32 graph = [[1, 2], [2, 3], [3, 4], []]
33
34 start_vertex = 0
35 distances = bfs_shortest_path(graph, start_vertex)
36 print(distances) # Output: [0, 1, 1, 2, 3]
37
38 # Graph taken from above
39 #Graph as adjacency list
40 #{'A': ['B', 'E'], 'B': ['A', 'C', 'D', 'E'], 'C': ['B', 'D'], 'D': ['B', 'C', 'E'], 'E': ['A', 'B', 'D']}
41 # Graph nodes remapped to integers
42 # A 0
43 # B 1
44 # C 2
45 # D 3
46 # E 4
47
48 graph = [[1, 4], [0, 2, 3, 4], [2, 3], [1, 2, 4], [0, 1, 3]]
49
50 start_vertex = 0
51 distances = bfs_shortest_path(graph, start_vertex)
52 print(distances) #
53 print(f"Distance from node (0) to node (4) is {distances[4]}")
[0, 1, 2, 3, 4]
[0, 1, 2, 3, 2]
Distance from node (0) to node (4) is 2
```

https://s32nbxekqbfbh7bnyimrvmr3i-dot-us-central1.dataproc.googleusercontent.com/gateway/d

OneLogin Final Project Overview and Test Course Groups: 2022-0822 D

Mouse Select Text Draw Stamp Spotlight Eraser Format Undo Redo Clear

Who can see what you share here? Recording On ass Neu

v261-student-348823 > w261

File Edit View Run Kernel Git Tabs Settings Help

DEMO9\_WORKBOOK\_MASTER.ipynb

1. MASTER Unit 9 - Graph Algorithms at Scale

1.0.1. Notebook Set-Up

2. Exercise 1. Graphs Overview

3. Warm-up questions:

4. Exercise 2. Data Structures Review.

5. Exercise 3. Graph Traversal.

5.1. Depth-First Search

5.2. Breadth-First Search

5.3. Other great algorithm books:

5.3.1. BFS

5.3.2. DFS

5.4. BFS versus DFS (no animation)

6. TASK: Write BFS Search Function In Apache Spark And Pandas Dataframe:

6.1. TODO: iterative search over undirected graph

7. Exercise 4. SSSP for unweighted graphs.

8. Exercise 5: Dijkstra's algorithm (SSSP for weighted graphs)

8.0.0.1. What about longest paths?

8.1. 8.1. 1

9. Exercise 6. Distributed SSSP

9.0.0.1. Graph algorithms typically involve:

9.0.0.2. Key questions:

9.0.1. Distributed SSSP Algorithm

OLD\_demo8\_OLD\_workbook.ipynb demo8\_workbook\_MASTEF.ipynb demo7\_workbook\_MASTER.ipynb hw2\_Workbook\_final\_master.ipynb hw4\_Workbook\_MASTER.ipynb

Explained [ ]

sample\_docs\_loc = f'{M03\_FOLDER}/sample\_docs.txt'  
!echo "(sample\_doc)" | gsutil cp - {sample\_docs\_loc}

Q = [1, 2, exp]

def bfs\_shortest\_path(graph, start, goal):  
 # keep track of explored nodes  
 explored = []  
 # keep track of all the paths to be checked  
 queue = deque()

See notebook for two alternative solutions to SSSP for unweighted graphs

if node not in explored:  
 neighbours = graph[node]  
 # go through all neighbour nodes, construct new path and  
 # push it into the queue  
 for neighbour in neighbours:  
 new\_path = list(path)  
 new\_path.append(neighbour)  
 queue.append(new\_path)  
 # return path if neighbour is goal  
 if neighbour == goal:  
 distance = len(new\_path)-1 # for an unweighted graph, distance is just the number of hops  
 return new\_path, distance  
  
 # mark node as explored  
 explored.append(node)  
  
 # in case there's no path between the 2 nodes  
 return "So sorry, but a connecting path doesn't exist :("

map reduce

Broo

# Use BFS to calculate shortest path for unweighted graphs

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** w261-student-348823 > w261
- File Menu:** File, Edit, View, Run, Kernel, Git, Tabs, Settings, Help
- File Tree:** DEMO9\_WORKBOOK\_MASTER.ipynb  
1. MASTER Unit 9 - Graph Algorithms at Scale  
    1.0.1. Notebook Set-Up  
    2. Exercise 1. Graphs Overview  
    3. Warm-up questions:  
    4. Exercise 2. Data Structures Review.  
    5. Exercise 3. Graph Traversal.  
        5.1. Depth-First Search  
        5.2. Breadth-First Search  
        5.3. Other great algorithm books:  
            5.3.1. BFS  
            5.3.2. DFS  
        5.4. BFS versus DFS (no animation)  
    6. TASK: Write BFS Search Function In Apache Spark And Pandas Dataframe:  
        6.1. TODO: iterative search over undirected graph  
    7. Exercise 4. SSSP for unweighted graphs.  
    8. Exercise 5: Dijkstra's algorithm (SSSP for weighted graphs)  
        8.0.0.1. What about longest paths?  
        8.1. 8.1. ¶  
    9. Exercise 6. Distributed SSSP  
        9.0.0.1. Graph algorithms typically involve:  
        9.0.0.2. Key questions:  
    9.0.1. Distributed SSSP Algorithm  
        9.0.1.1. Node STATES
- Code Editor:** demo9\_workbook\_MASTERF ●  
The code in the editor is a Python function for Breadth-First Search (BFS) to find the shortest path in an unweighted graph. It uses a queue to store paths and explores nodes until it finds the goal or exhausts all possible paths.
- Python 3 Environment:** Python 3 | Idle
- Status Bar:** Mode: Edit | In 44 Col 47 demo9 workbook MASTER

```
def bfs_shortest_path(graph, start, goal):
    # keep track of explored nodes
    explored = []
    # keep track of all the paths to be checked
    queue = deque()
    queue.append([start])

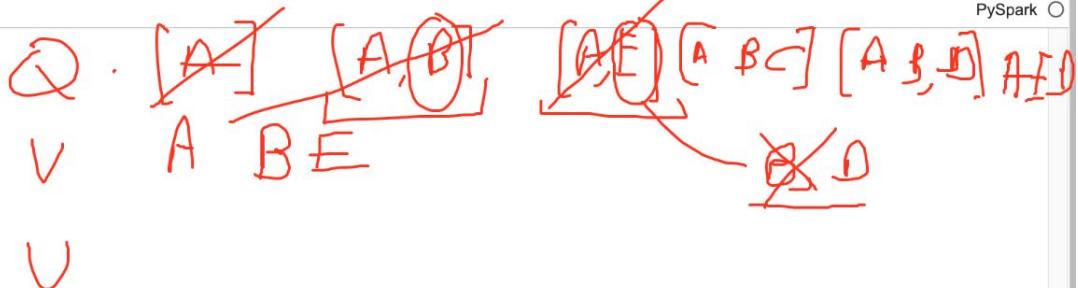
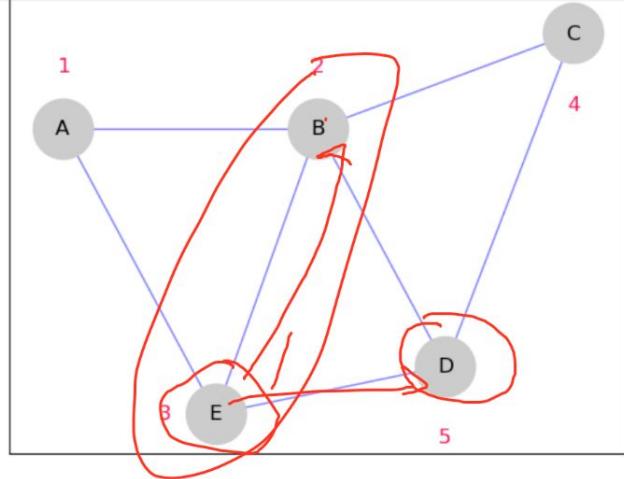
    distance = 0
    # return path if start is goal
    if start == goal:
        return "That was easy! Start = goal"

    # keeps looping until all possible paths have been checked
    while queue: # que is a list of paths where the last element is the node of interest
        # pop the first path from the queue
        path = queue.popleft()
        # get the last node from the path
        node = path[-1]
        if node not in explored:
            neighbours = graph[node] #neighbors
            # go through all neighbour nodes, construct a new path and
            # push it into the queue
            for neighbour in neighbours:
                new_path = list(path)
                new_path.append(neighbour)
                # que is a list of paths where the last element is the node of interest
                # queue consists of [A B]; the node in the queue is B
                queue.append(new_path)
            # return path if neighbour is goal
            if neighbour == goal:
                distance = len(new_path)-1 # for an unweighted graph, distance is just the number of hops
                return new_path, distance

        # mark node as explored
        explored.append(node)

    # in case there's no path between the 2 nodes
    return "So sorry, but a connecting path doesn't exist :("

print("BFS Shortest paths")
print(bfs_shortest_path(ADJ_GRAPH, 'A', 'D'))
```



```
[20]: 1 ADJ_GRAPH
2 #Graph as adjacency list
3 #{'A': ['B', 'E'], 'B': ['A', 'C', 'D', 'E'], 'C': ['B', 'D'], 'D': ['B', 'C', 'E'], 'E': ['A', 'B', 'D']}
4 ADJ_GRAPH = {'A': ['B', 'E'],
5   'B': ['A', 'C', 'D', 'E'],
6   'C': ['B', 'D'],
7   'D': ['B', 'C', 'E'],
8   'E': ['A', 'B', 'D']}
9 ADJ_GRAPH
```

```
[20]: {'A': ['B', 'E'],
'B': ['A', 'C', 'D', 'E'],
'C': ['B', 'D'],
'D': ['B', 'C', 'E'],
'E': ['A', 'B', 'D']}
```

```
[ ]: 1
```

```
[27]: 1 # finds shortest path between 2 nodes of a graph using BFS
2 # Does not consider edge weights
3
4 def bfs_shortest_path(graph, start, goal):
5   .....
6
7   # keep track of explored nodes
8   explored = []
```

s Settings Help

OLD\_demo8\_OLD\_workbo X demo8\_workbook\_MASTEF X demo7\_workbook\_MASTER X hw3\_Workbook\_final\_mast X hw4\_Workbook\_MASTER.iX

ale + % C &gt; Code git

/ [1 files] 0.0 B/ 0.0 B  
Operation completed over 1 objects.

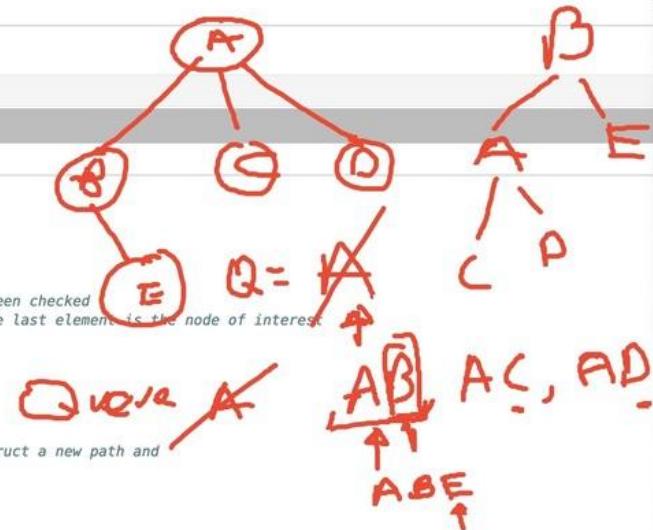
[22]:

demo9\_workbook\_MASTEF X c02w06\_BFS\_Spark\_SQL.iX

ale + % C &gt; Code git

```
10
11     distance = 0
12     # return path if start is goal
13     if start == goal:
14         return "That was easy! Start = goal"
15
16     # keeps looping until all possible paths have been checked
17     while queue: # que is a list of paths where the last element is the node of interest
18         # pop the first path from the queue
19         path = queue.pop(0)
20         # get the last node from the path
21         node = path[-1]
22         if node not in explored:
23             neighbours = graph[node] #neighbors
24             # go through all neighbour nodes, construct a new path and
25             # push it into the queue
26             for neighbour in neighbours:
27                 new_path = list(path)
28                 new_path.append(neighbour)
29                 # que is a list of paths where the last element is the node of interest
30                 # queue consists of [A B]; the node in the queue is B
31                 queue.append(new_path)
32                 # return path if neighbour is goal
33                 if neighbour == goal:
34                     distance = len(new_path)-1 # for an unweighted graph, distance is just the number of hops
35                     return new_path, distance
36
37             # mark node as explored
38             explored.append(node)
39
40     # in case there's no path between the 2 nodes
41     return "So sorry, but a connecting path doesn't exist :("
42
43 print("BFS Shortest paths")
44 print(bfs_shortest_path(ADJ_GRAPH, 'A', 'D'))
```

BFS Shortest paths



## Possible solution 1 to BFS/SSSP/Dijkstra with a single Spark Dataframe

Implement BFS at scale to calculate shortest path for unweighted graphs

queue\_len = 0

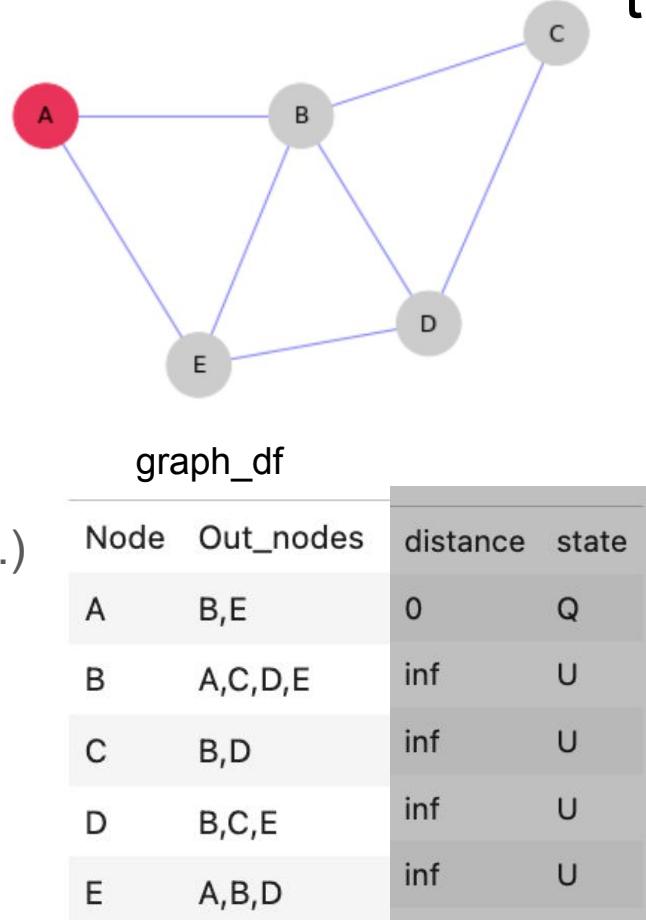
graph\_df = load\_graph(location)

status\_df =

While queue\_len > 0:

graph....(status\_df).flatmap(....)

reduce



# Iterate through rows in a dataframe: map(), mapPartitions()

PySpark provides map(), mapPartitions() to loop/iterate through rows in RDD/DataFrame to perform the complex transformations, and these two returns the same number of records as in the original DataFrame but the number of columns could be different (after add/update).

PySpark also provides foreach() & foreachPartitions() actions to loop/iterate through each Row in a DataFrame but these two returns nothing.

This article explains how to use these methods to get DataFrame column values and process them:

<https://sparkbyexamples.com/pyspark/pyspark-loop-iterate-through-rows-in-dataframe/>

# DataFrames: mini review: create a dataframe

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName('SparkByExamples.com').getOrCreate()

data = [('James','Smith','M',30), ('Anna','Rose','F',41),
        ('Robert','Williams','M',62),
       ]
columns = ["firstname","lastname","gender","salary"]
df = spark.createDataFrame(data=data, schema = columns)
df.show()
+-----+-----+-----+
|firstname|lastname|gender|salary|
+-----+-----+-----+
|    James|   Smith|     M|     30|
|    Anna|    Rose|     F|     41|
| Robert|Williams|     M|     62|
+-----+-----+-----+
```

```
from pyspark.sql.functions import concat_ws,col,lit
df.select(concat_ws(",") , df.firstname,df.lastname).alias("name"), \
          Df.gender, lit(df.salary*2).alias("new_salary")).show()
+-----+-----+-----+
|      name|gender|new_salary|
+-----+-----+-----+
| James,Smith|     M|       60|
| Anna,Rose|     F|       82|
|Robert,Williams|     M|      124|
```

Mostly for simple computations, instead of iterating through using map() and foreach(), you should use either [DataFrame select\(\)](#) or [DataFrame withColumn\(\)](#) in conjunction with PySpark SQL functions.

# map() to loop/iterate through the PySpark DataFrame/RDD

[PySpark map\(\) Transformation](#) is used to loop/iterate through the PySpark DataFrame/RDD by applying the transformation function (lambda) on every element (Rows and Columns) of RDD/DataFrame. PySpark doesn't have a map() in DataFrame instead it's in RDD hence we need to convert DataFrame to RDD first and then use the map(). It returns an RDD and you should [Convert RDD to PySpark DataFrame](#) if needed.

If you have a heavy initialization use PySpark mapPartitions() transformation instead of map(), as with mapPartitions() heavy initialization executes only once for each partition instead of every record.

```
# Referring columns by index.  
rdd=df.rdd.map(lambda x: (x[0]+", "+x[1],x[2],x[3]*2))  
df2=rdd.toDF(["name","gender","new_salary"])  
df2.show()  
# Referring Column Names  
rdd2=df.rdd.map(lambda x: (x.firstname+" "+x.lastname,x.gender,x.salary*2) )  
  
# By Calling function  
def func1(x):  
    name=x.firstname+" "+x.lastName  
    return (name, x.gender.lower(), x.salary*2)  
  
rdd2=df.rdd.map(lambda x: func1(x))
```

```

data = [('James','Smith','M',30),
        ('Anna','Rose','F',41),
        ('Robert','Williams','M',62),
        ]

columns = ["firstname","lastname","gender","salary"]
df = spark.createDataFrame(data=data, schema = columns)
df.show()
+-----+-----+-----+
|firstname|lastname|gender|salary|
+-----+-----+-----+
|      James|    Smith|      M|     30|
|      Anna|     Rose|      F|     41|
|   Robert|Williams|      M|     62|
+-----+-----+-----+

```

X is a row in the dataframe

```

# Refering columns by index.
rdd2=df.rdd.map(lambda x:
                 (x[0]+", "+x[1],x[2],x[3]*2)
                )
df2=rdd2.toDF(["name","gender","new_salary"])
df2.show()
+-----+-----+-----+
|         name|gender|new_salary|
+-----+-----+-----+
| James,Smith|      M|       60|
|   Anna,Rose|      F|       82|
|Robert,Williams|      M|      124|
+-----+-----+-----+

```

```

data = [('James','Smith','M',30),
 ('Anna','Rose','F',41),
 ('Robert','Williams','M',62),
]

columns = ["firstname","lastname","gender","salary"]
df = spark.createDataFrame(data=data, schema = columns)
df.show()
+-----+-----+-----+
|firstname|lastname|gender|salary|
+-----+-----+-----+
|    James|   Smith|     M|    30|
|    Anna|    Rose|     F|    41|
| Robert|Williams|     M|    62|
+-----+-----+-----+

```

```

# Referring columns by index.
rdd2=df.rdd.map(lambda x:
    (x[0]+", "+x[1],x[2],x[3]*2)
)
df2=rdd2.toDF(["name","gender","new_salary"])
df2.show()
+-----+-----+-----+
|         name|gender|new_salary|
+-----+-----+-----+
| James,Smith|     M|        60|
|   Anna,Rose|     F|        82|
| Robert,Will|     M|       124|
+-----+-----+-----+

```

Note that above we used the index to get the column values, alternatively, you can also refer to the DataFrame column names while iterating.

Copy

```

# Referring Column Names
rdd2=df.rdd.map(lambda x:
    (x["firstname"]+", "+x["lastname"],x["gender"],x["salary"]*2)
)

```

Another alternative

```

# Referring Column Names
rdd2=df.rdd.map(lambda x:
    (x.firstname+", "+x.lastname,x.gender,x.salary*2)
)

```

You can also create a custom function to perform an operation. Below `func1()` function executes for every DataFrame row from the lambda function.

[https://sparkbyexamples.com/pyspark/pyspark-map-transformation/?expand\\_article=1](https://sparkbyexamples.com/pyspark/pyspark-map-transformation/?expand_article=1)

# df.foreach

```
# Foreach example
def f(x): print(x)
df.foreach(f)

# Another example
df.foreach(lambda x:
    print("Data==>" +x["firstname"] +", "+x["lastname"] +", "+x["gender"] +", "+str(x["salary"] *2))
)
```

## Possible solution 1 to BFS/SSSP/Dijkstra with a single Spark Dataframe

Implement BFS at scale to calculate shortest path for unweighted graphs

queue\_len = 0

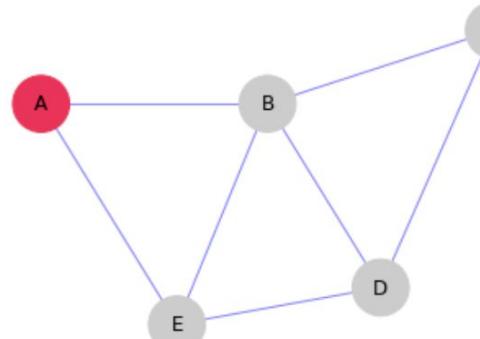
graph\_df = load\_graph(location)

status\_df =

While queue\_len > 0:

graph.....(status\_df).flatmap(.....)

reduce



graph_df				
Node	Out_nodes	distance	state	
A	B,E	0	Q	
B	A,C,D,E	inf	U	
C	B,D	inf	U	
D	B,C,E	inf	U	
E	A,B,D	inf	U	

## Possible solution 2 to BFS/SSSP/Dijkstra with a single Spark DataframeS

Implement BFS at scale to calculate shortest path for unweighted graphs

queue\_len = 0

graph\_df = load\_graph(location)

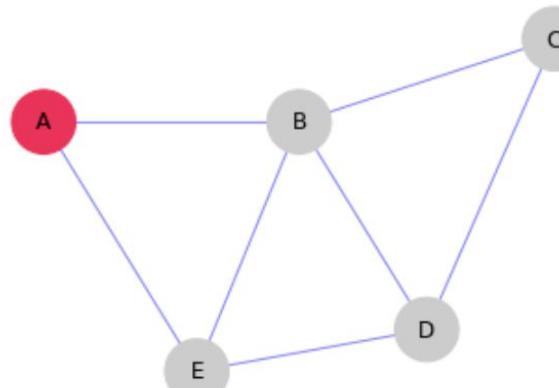
status\_df =

While queue\_len > 0:

graph.....(status\_df).flatmap(.....)

reduce

JOIN



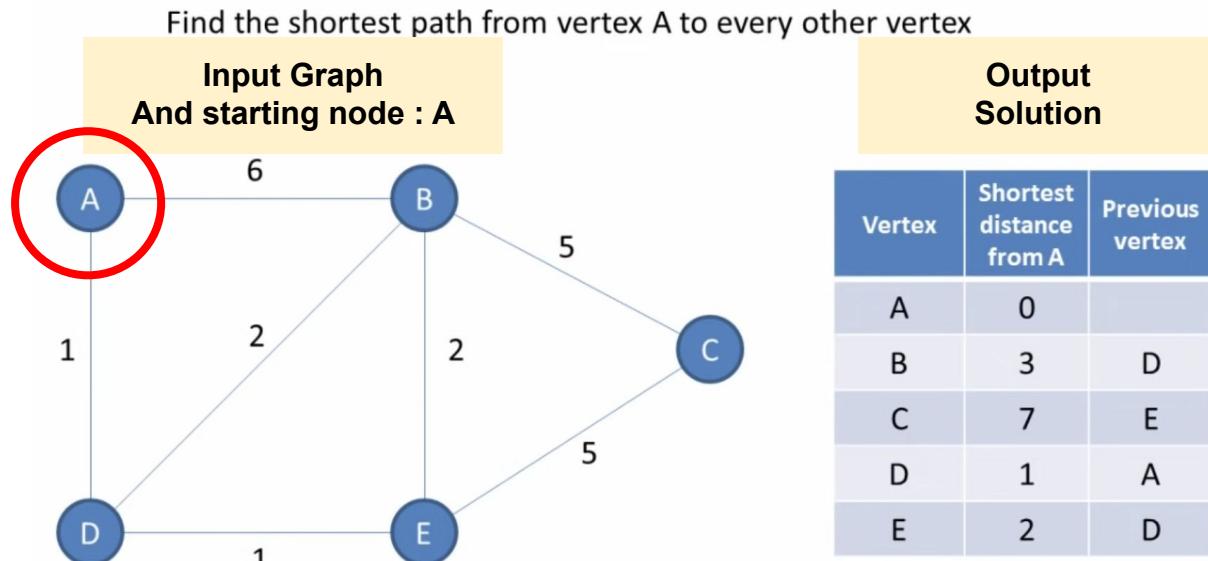
graph_df		status_df		
Node	Out_nodes	Node	distance	state
A	B,E	A	0	Q
B	A,C,D,E	B	inf	U
C	B,D	C	inf	U
D	B,C,E	D	inf	U
E	A,B,D	E	inf	U

---

# Dijkstra's shortest path algorithm for a weighted graph

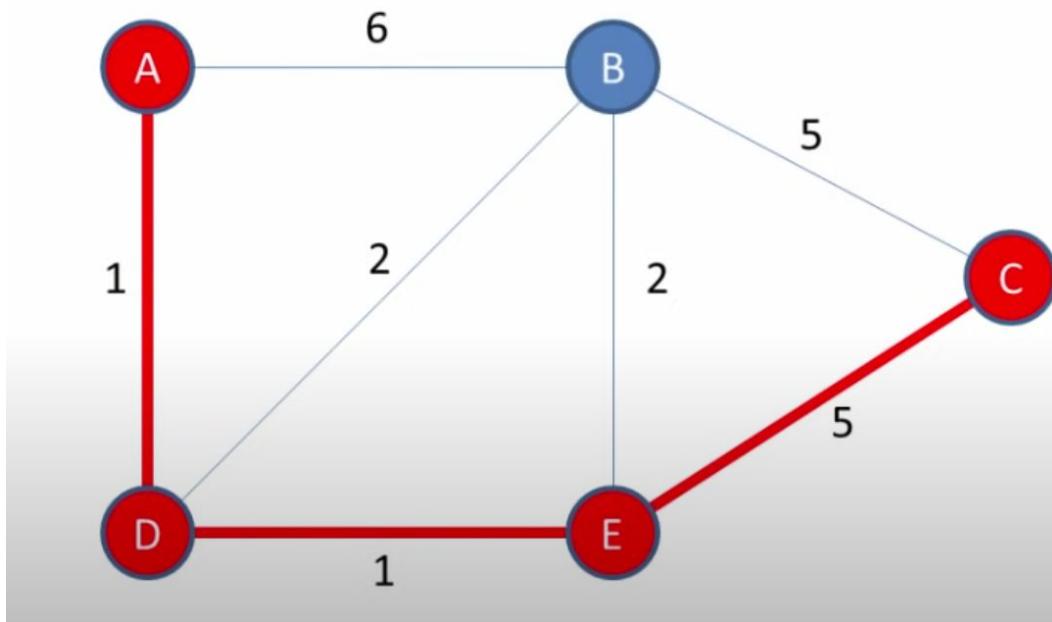
# Dijkstra's shortest path algorithm for a weighted graph

- Dijkstra's shortest path algorithm generates a set of information that includes the shortest paths from a starting vertex and every other vertex in the graph.
- Why Dijkstra's shortest path algorithm is an example of a greedy algorithm.



<https://www.youtube.com/watch?v=pVfj6mxhdMw>

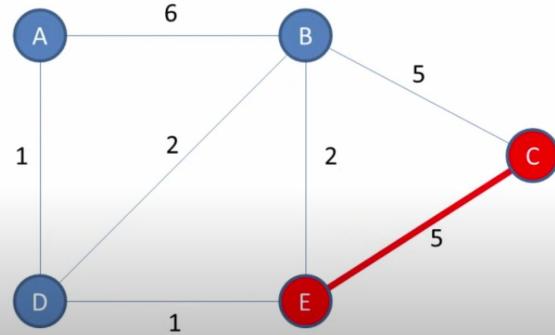
A to C



Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

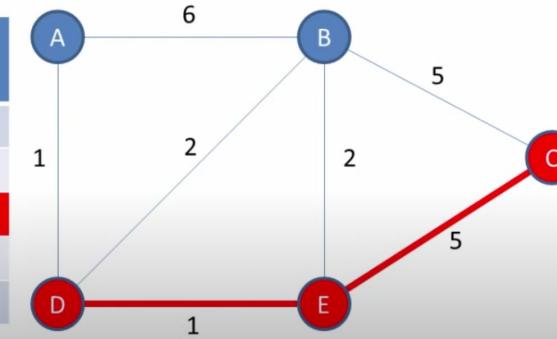
We arrived at C via E

Path = E → C



We arrived at E via D

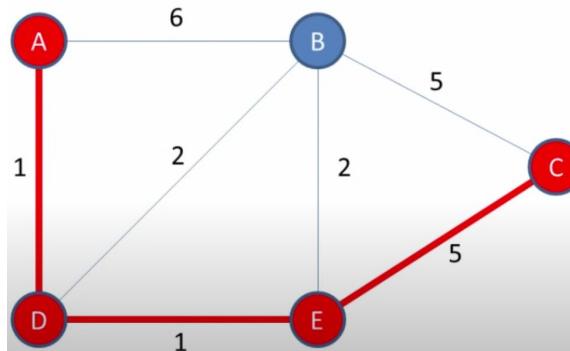
Path = D → E → C



Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

We arrived at D via A

Path = A → D → E → C



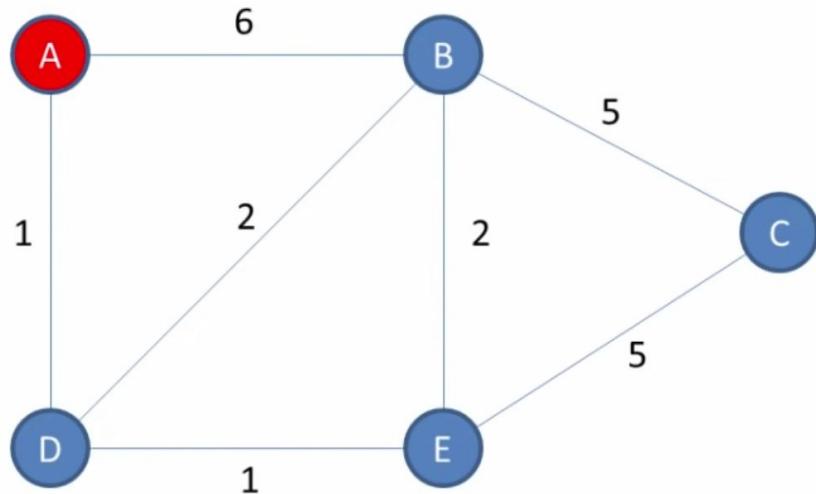
Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

Consider the start vertex, A

Distance to A from A = 0

Distances to all other vertices from A are unknown, therefore  $\infty$  (infinity)

A node's State can be  
{Queue, Visited, Unvisited}

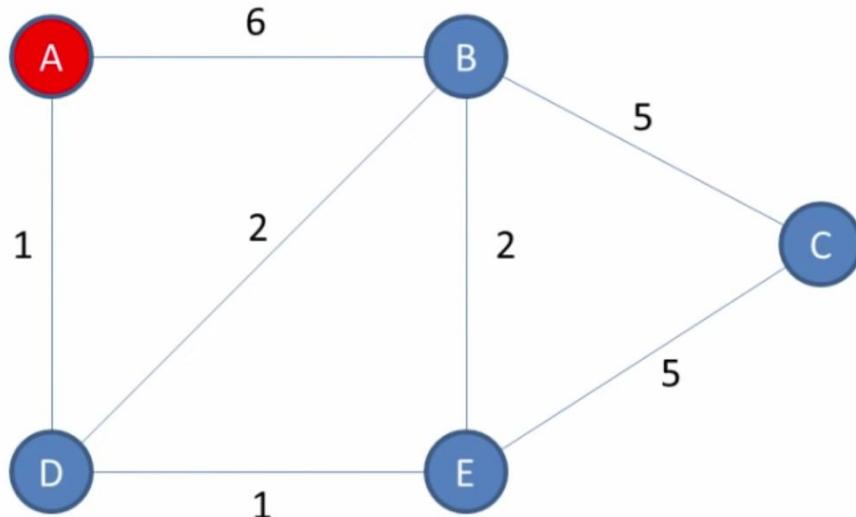


Vertex	Shortest distance from A	Previous vertex	State
A			Q
B			U
C			U
D			U
E			U

Visited = [ ]

Unvisited = [A, B, C, D, E]

Visit the unvisited vertex with the smallest known distance from the start vertex  
*First time around, this is the start vertex itself, A*

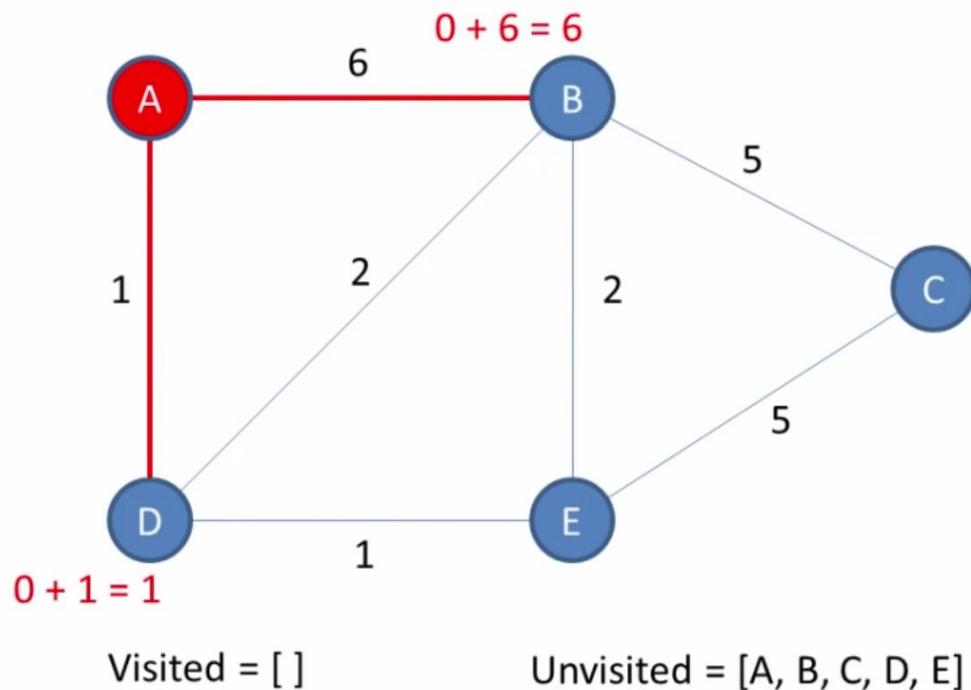


Vertex	Shortest distance from A	Previous vertex	State
A	0		Q
B	$\infty$		U
C	$\infty$		U
D	$\infty$		U
E	$\infty$		U

Visited = [ ]

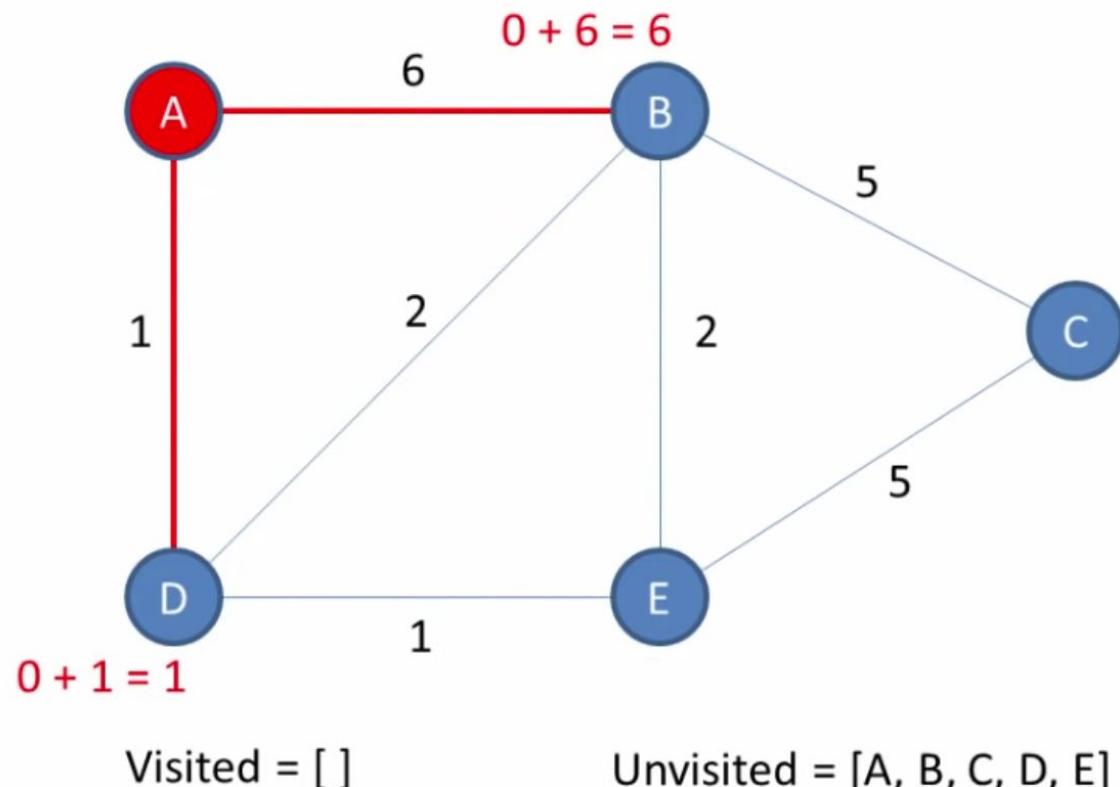
Unvisited = [A, B, C, D, E]

For the current vertex, calculate the distance of each neighbour from the start vertex



Vertex	Shortest distance from A	Previous vertex
A	0	
B	$\infty$	
C	$\infty$	
D	$\infty$	
E	$\infty$	

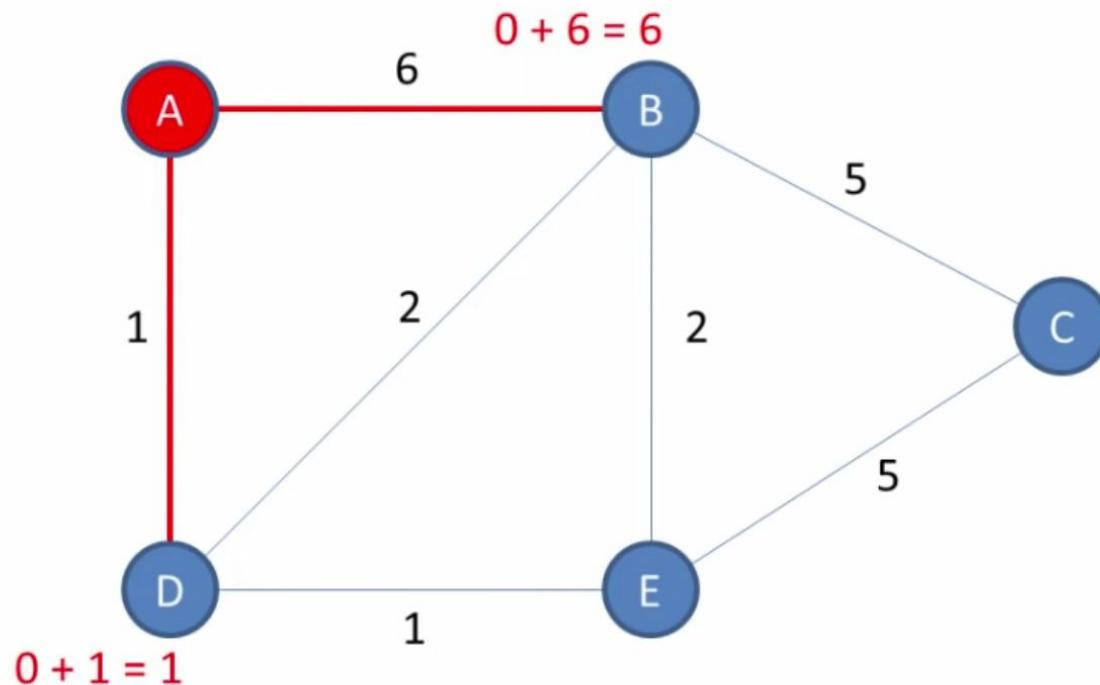
If the calculated distance of a vertex is less than the known distance, update the shortest distance



Vertex	Shortest distance from A	Previous vertex
A	0	
B	6	
C	$\infty$	
D	1	
E	$\infty$	

Update the previous vertex for each of the updated distances

*In this case we visited B and D via A*

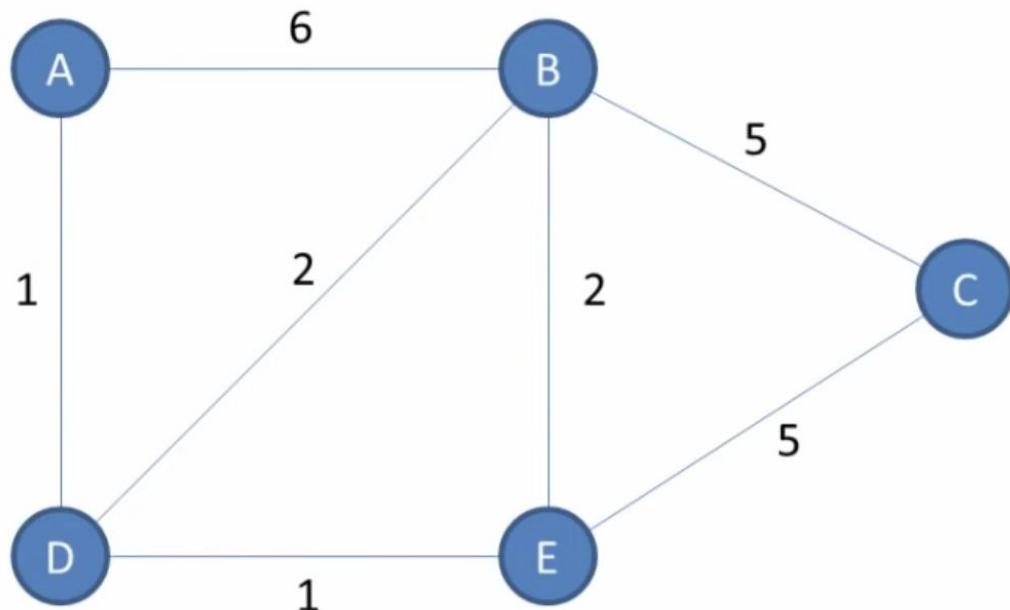


Visited = [ ]

Unvisited = [A, B, C, D, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	6	
C	$\infty$	
D	1	
E	$\infty$	

Add the current vertex to the list of visited vertices



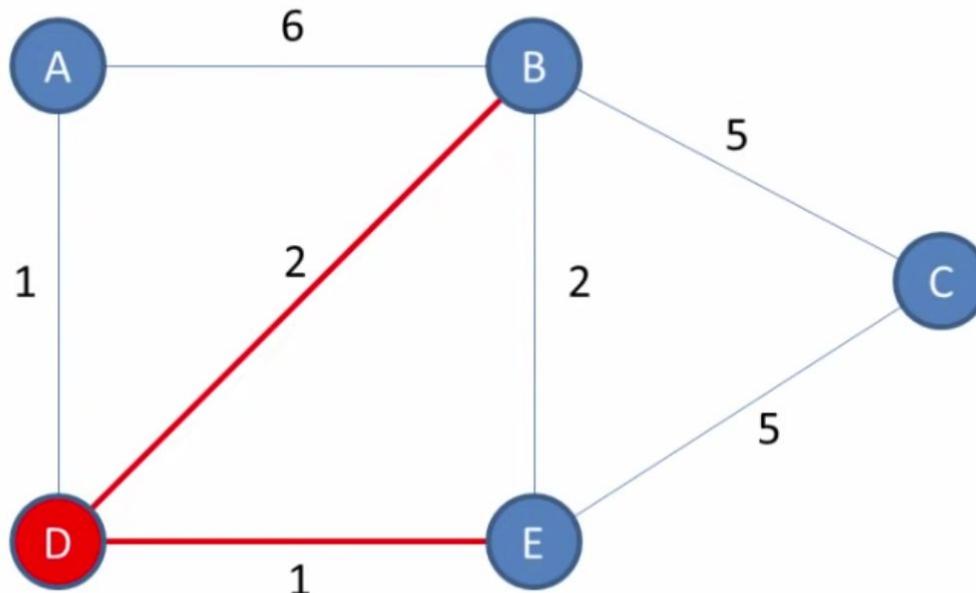
Visited = [A]

Unvisited = [B, C, D, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	6	A
C	$\infty$	
D	1	A
E	$\infty$	

For the current vertex, examine its unvisited neighbours

We are currently visiting D and its unvisited neighbours are B and E

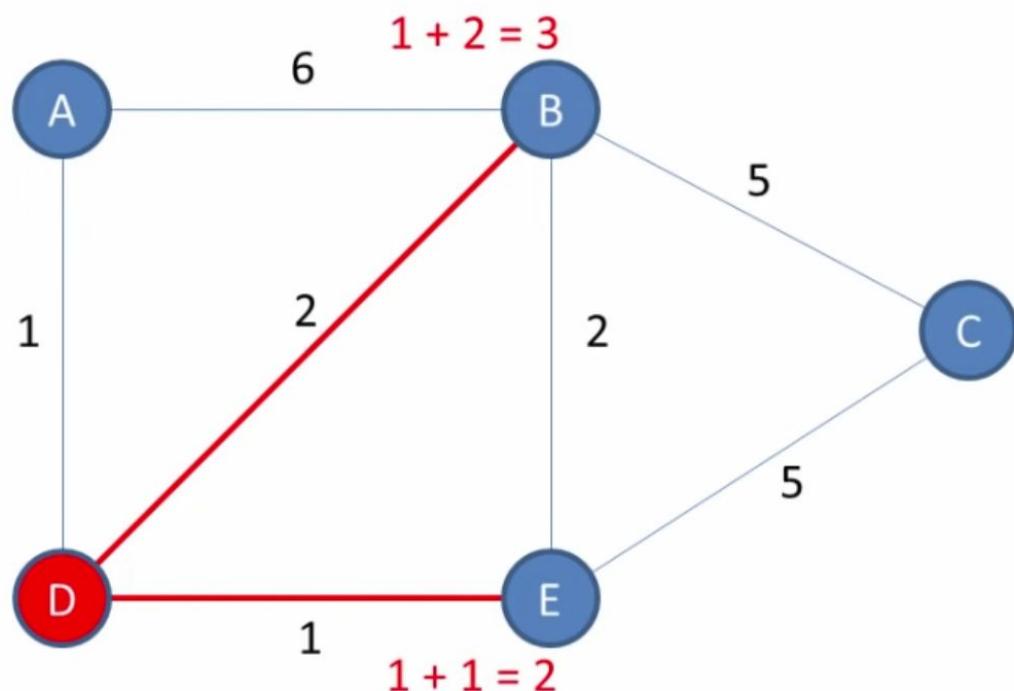


Vertex	Shortest distance from A	Previous vertex
A	0	
B	6	A
C	$\infty$	
D	1	A
E	$\infty$	

Visited = [A]

Unvisited = [B, C, D, E]

For the current vertex, calculate the distance of each neighbour from the start vertex

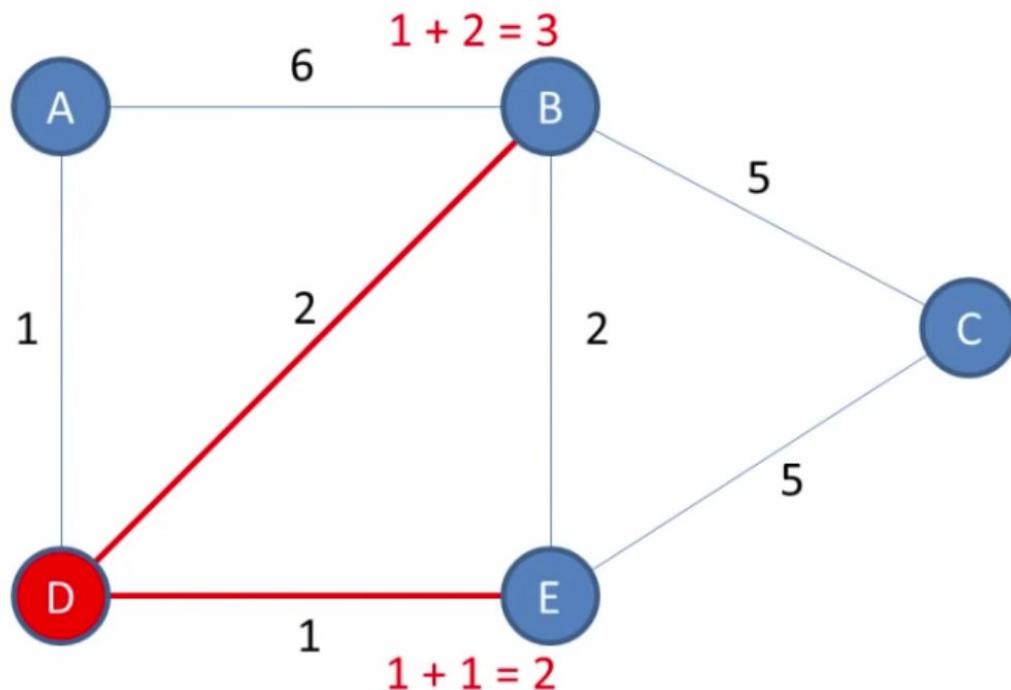


Vertex	Shortest distance from A	Previous vertex
A	0	
B	6	A
C	$\infty$	
D	1	A
E	$\infty$	

Visited = [A]

Unvisited = [B, C, D, E]

If the calculated distance of a vertex is less than the known distance, update the shortest distance

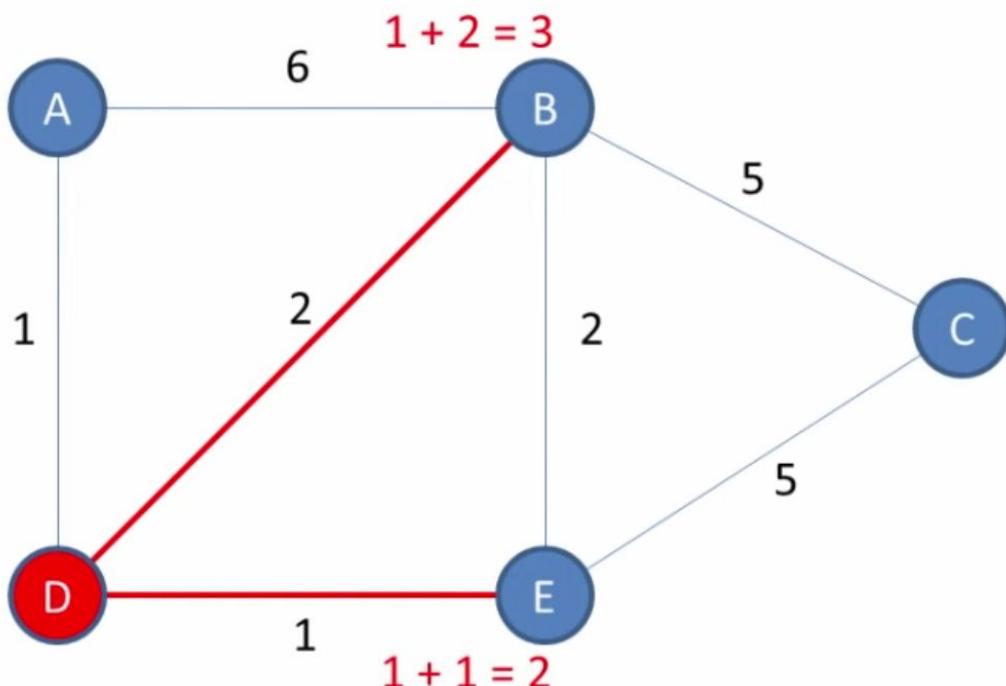


Visited = [A]

Unvisited = [B, C, D, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	6	A
C	$\infty$	
D	1	A
E	$\infty$	

If the calculated distance of a vertex is less than the known distance, update the shortest distance

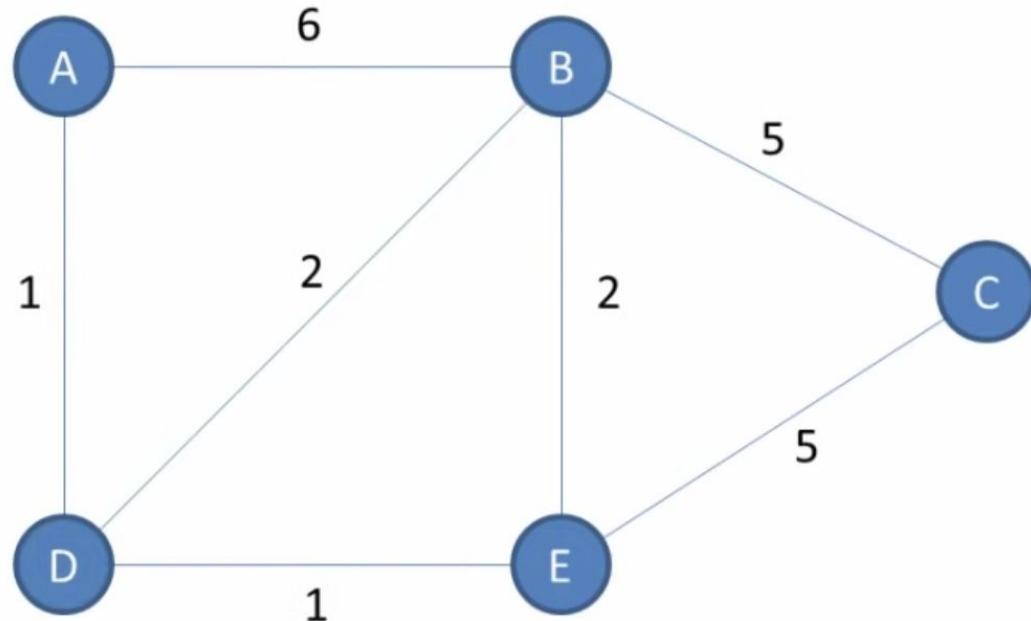


Visited = [A]

Unvisited = [B, C, D, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	A
C	$\infty$	
D	1	A
E	2	

Add the current vertex to the list of visited vertices

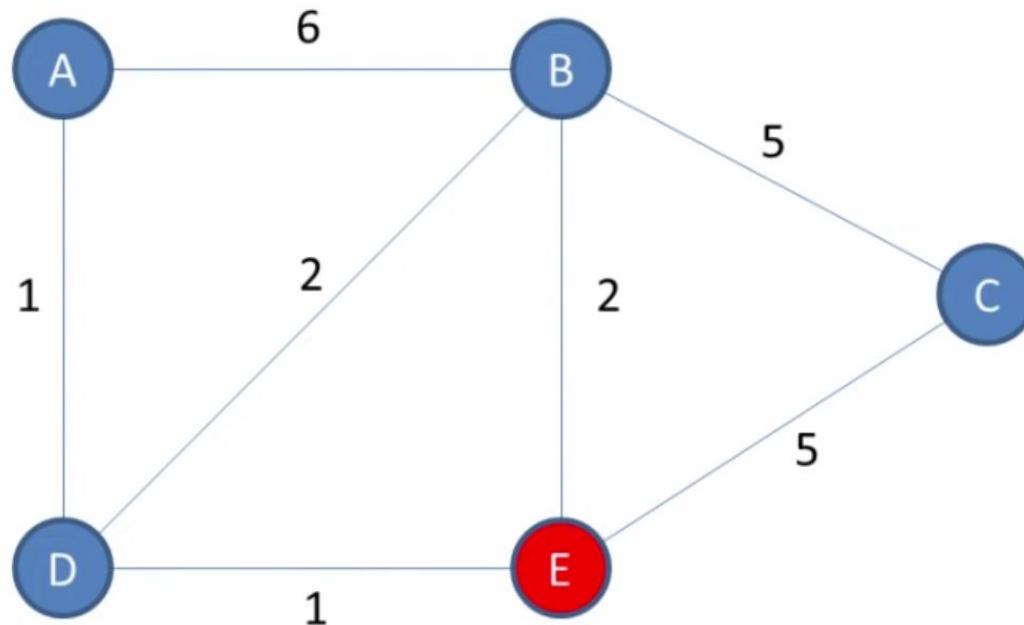


Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	$\infty$	
D	1	A
E	2	D

Visited = [A, D]

Unvisited = [B, C, E]

Visit the unvisited vertex with the smallest known distance from the start vertex  
*This time around, it is vertex E*



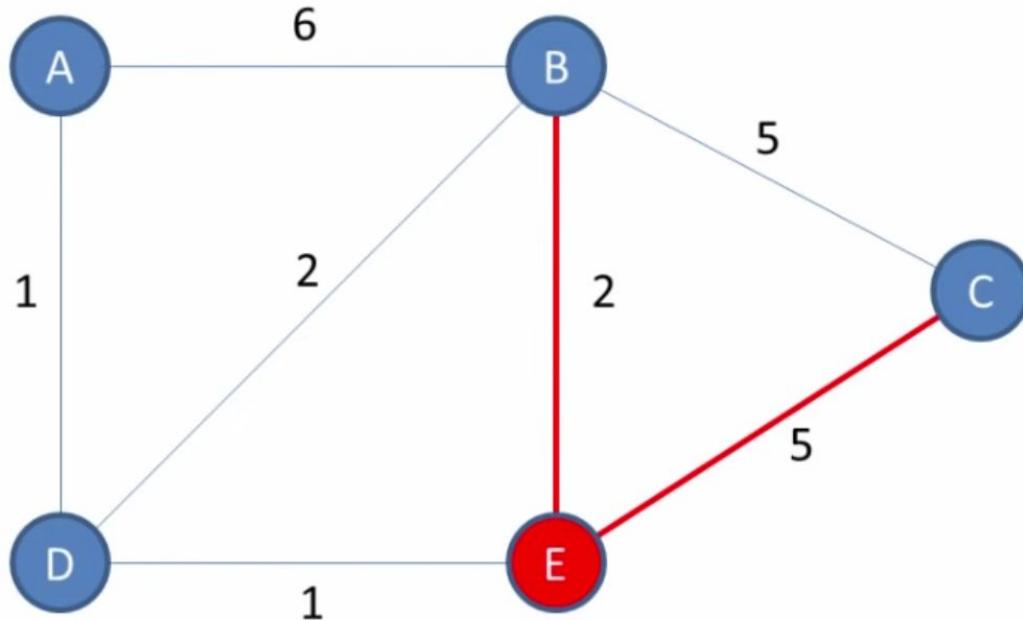
Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	$\infty$	
D	1	A
E	2	D

Visited = [A, D]

Unvisited = [B, C, E]

For the current vertex, examine its unvisited neighbours

We are currently visiting E and its unvisited neighbours are B and C

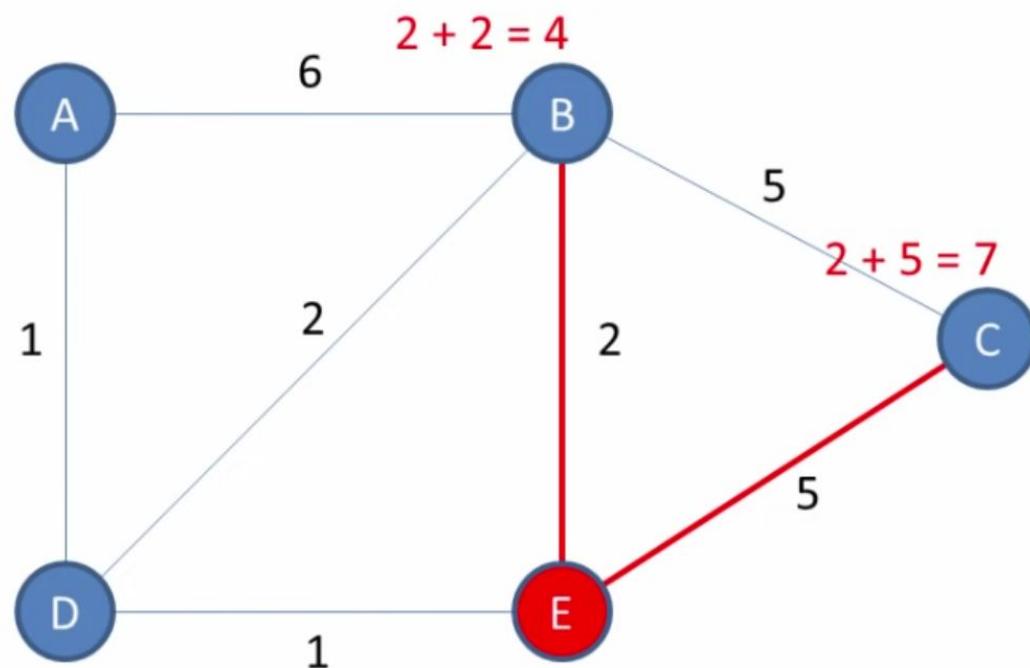


Visited = [A, D]

Unvisited = [B, C, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	$\infty$	
D	1	A
E	2	D

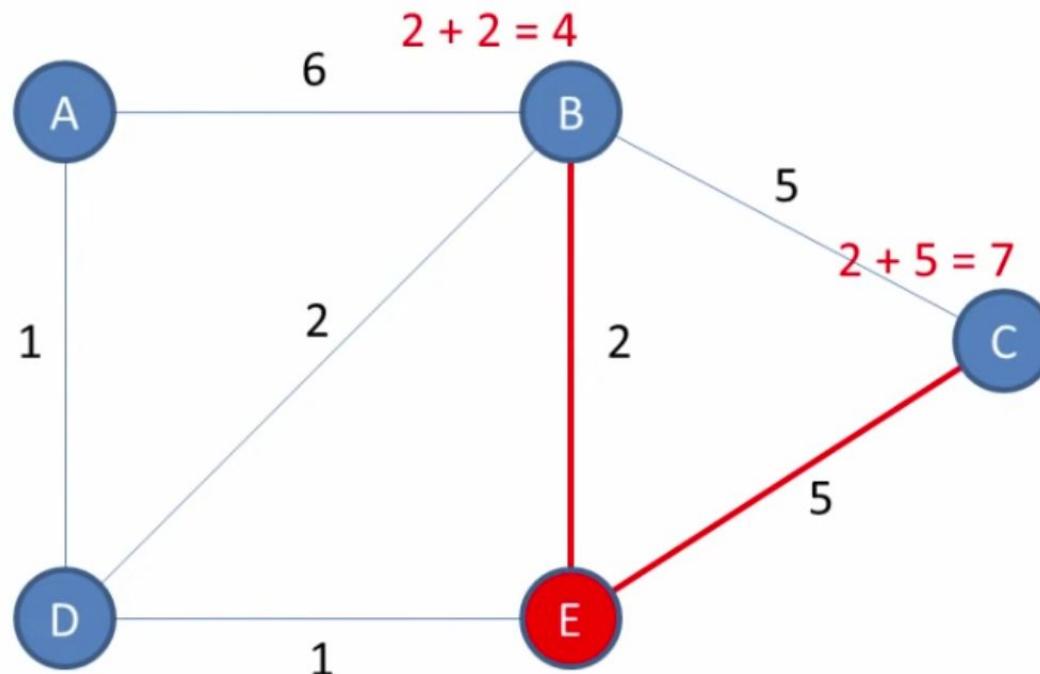
For the current vertex, calculate the distance of each neighbour from the start vertex



Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	$\infty$	
D	1	A
E	2	D

If the calculated distance of a vertex is less than the known distance, update the shortest distance

*We do not need to update the distance to B*



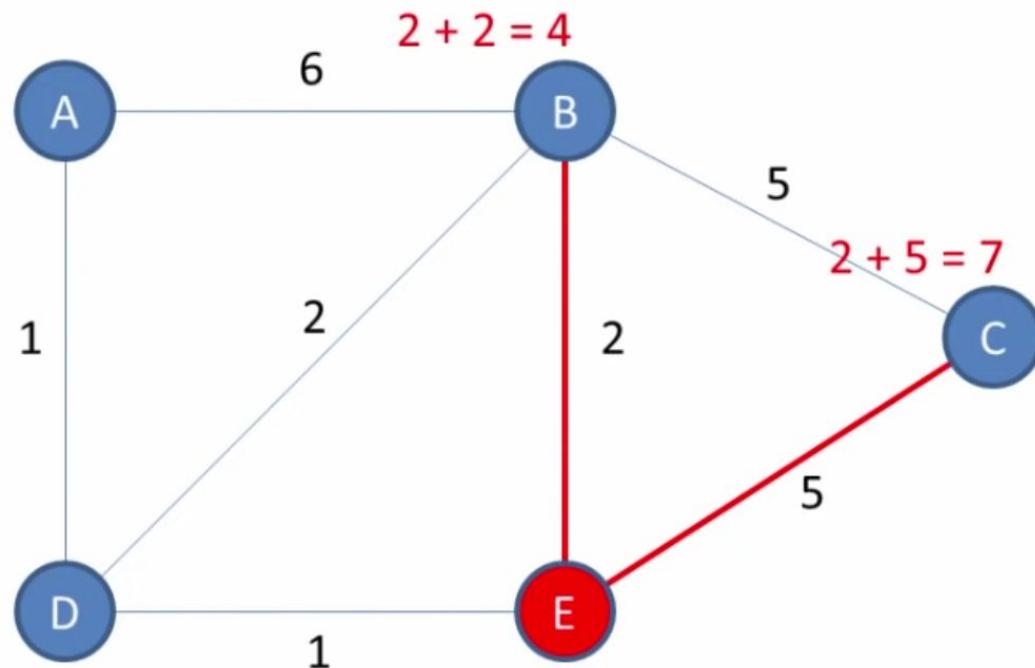
Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	$\infty$	
D	1	A
E	2	D

Visited = [A, D]

Unvisited = [B, C, E]

Update the previous vertex for each of the updated distances

In this case we visited C via E

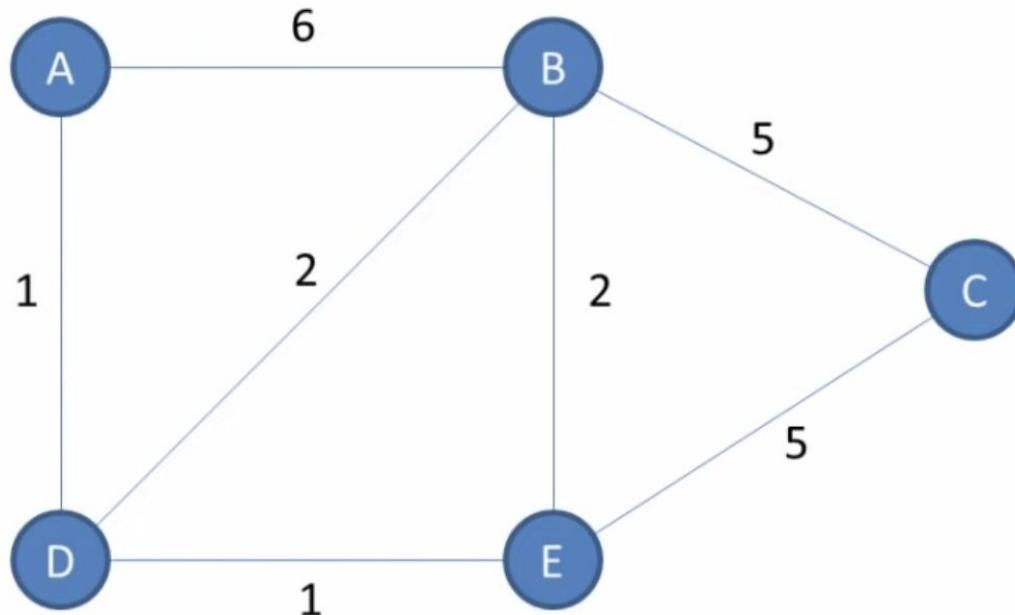


Visited = [A, D]

Unvisited = [B, C, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	
D	1	A
E	2	D

Add the current vertex to the list of visited vertices



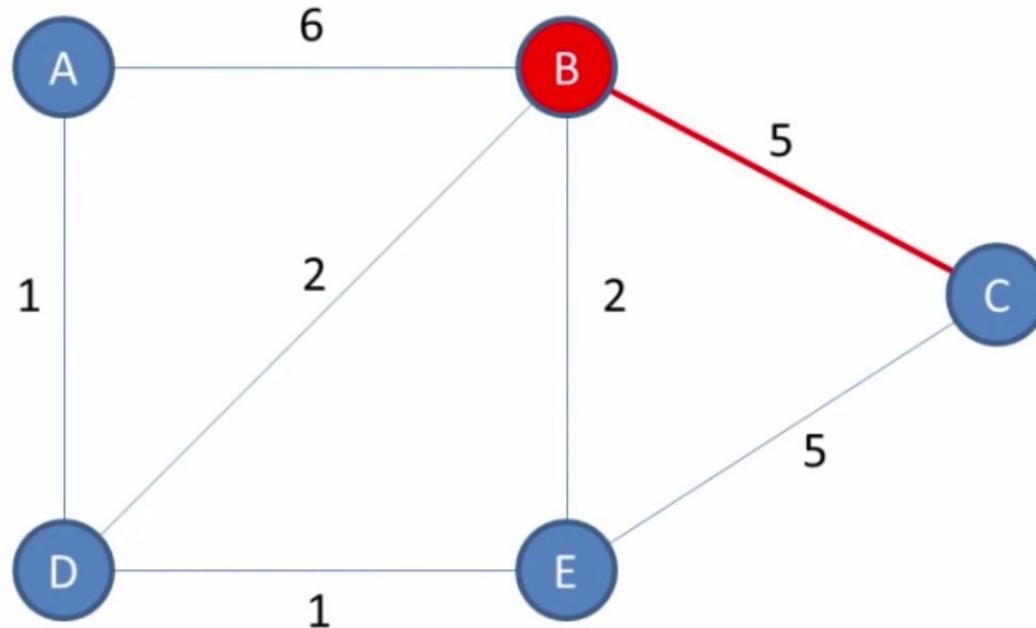
Visited = [A, D, E]

Unvisited = [B, C]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

For the current vertex, examine its unvisited neighbours

We are currently visiting B and its only unvisited neighbour is C

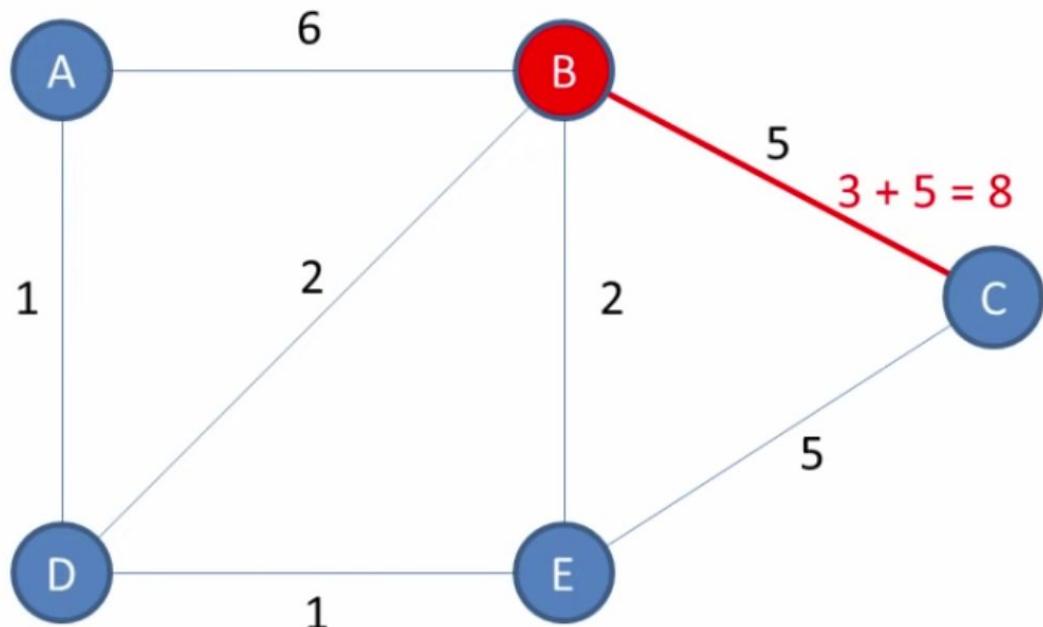


Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

Visited = [A, D, E]

Unvisited = [B, C]

For the current vertex, calculate the distance of each neighbour from the start vertex



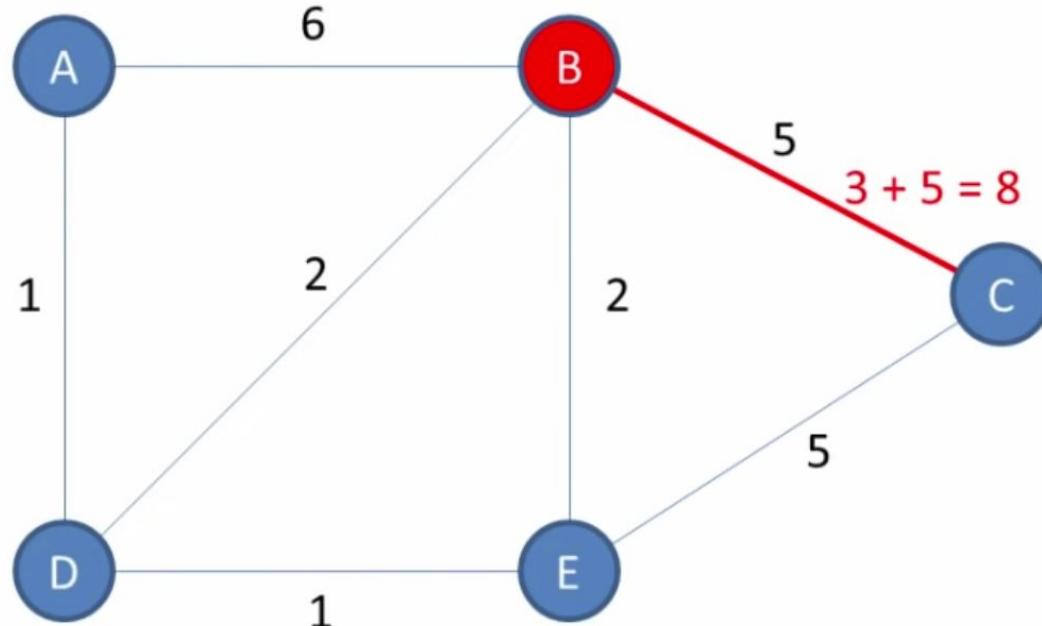
Visited = [A, D, E]

Unvisited = [B, C]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

If the calculated distance of a vertex is less than the known distance, update the shortest distance

We do not need to update the distance to C

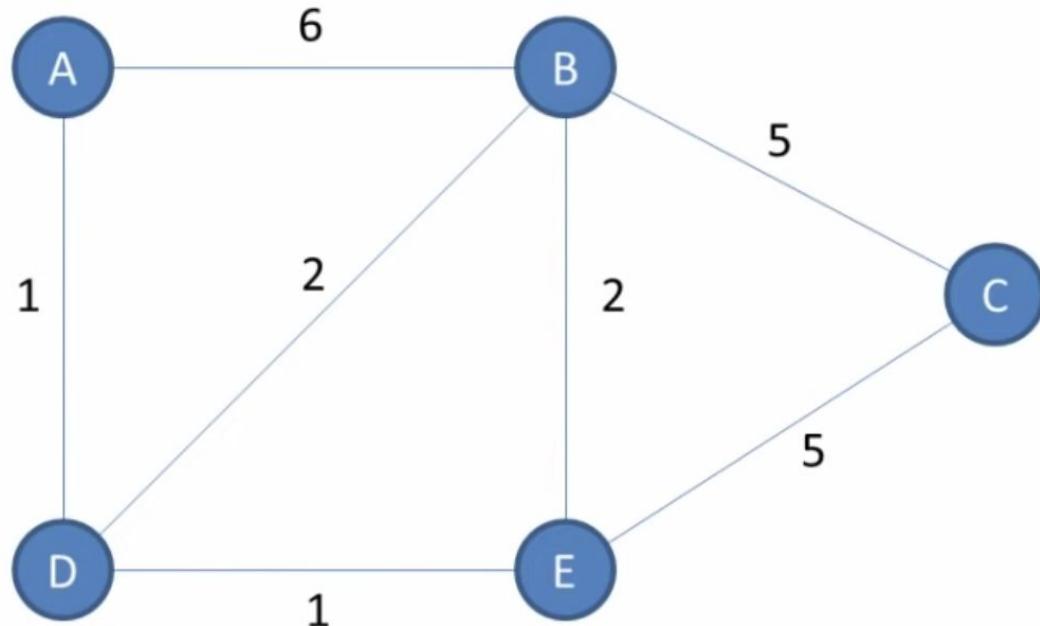


Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

Visited = [A, D, E]

Unvisited = [B, C]

Add the current vertex to the list of visited vertices



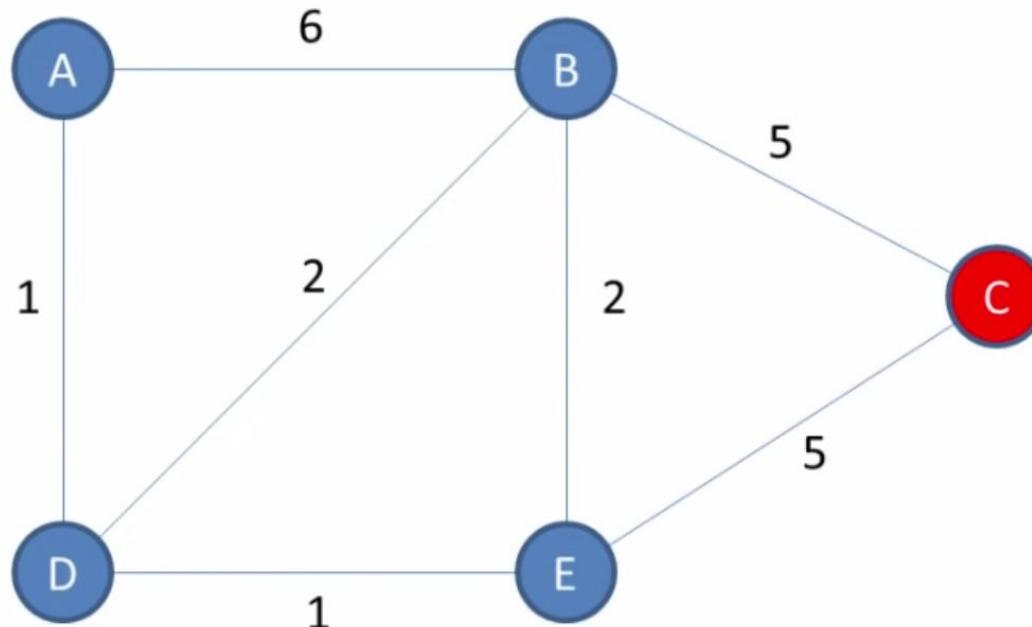
Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

Visited = [A, D, E, B]

Unvisited = [C]

For the current vertex, examine its unvisited neighbours

We are currently visiting C and it has no unvisited neighbours

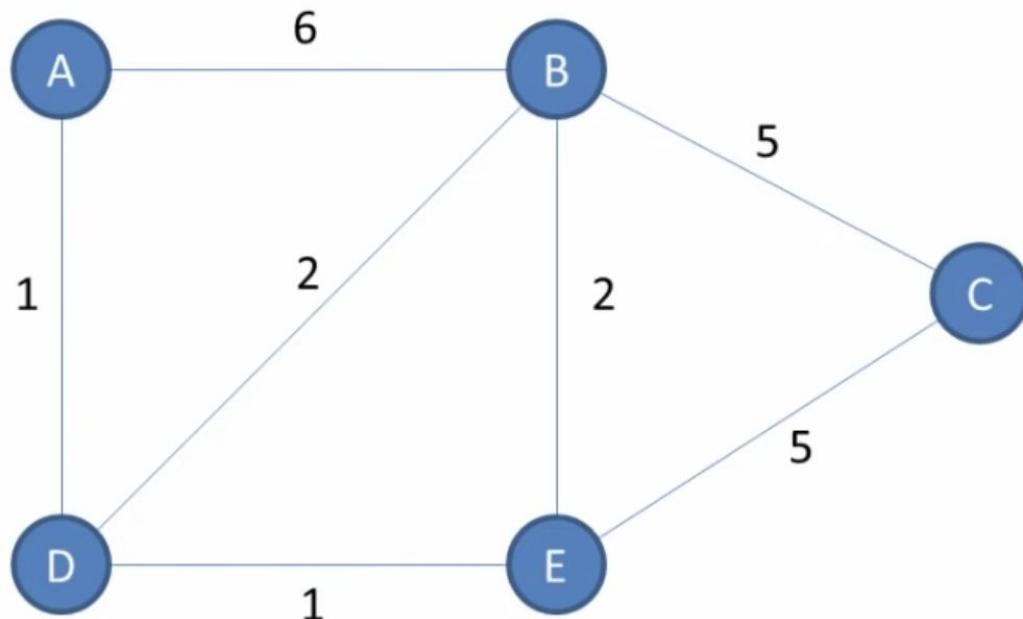


Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

Visited = [A, D, E, B]

Unvisited = [C]

Add the current vertex to the list of visited vertices



Vertex	Shortest distance from A	Previous vertex
A	0	
B	3	D
C	7	E
D	1	A
E	2	D

Visited = [A, D, E, B, **C**] Unvisited = []

# Dikstras Algo: greedy algo (make locally optimal decisions)

Let distance of start vertex from start vertex = 0

Let distance of all other vertices from start =  $\infty$  (infinity)

Repeat

Visit the unvisited vertex with the smallest known distance from the start vertex

For the current vertex, examine its unvisited neighbours

For the current vertex, calculate distance of each neighbour from start vertex

If the calculated distance of a vertex is less than the known distance, update the shortest distance

Update the previous vertex for each of the updated distances

Add the current vertex to the list of visited vertices

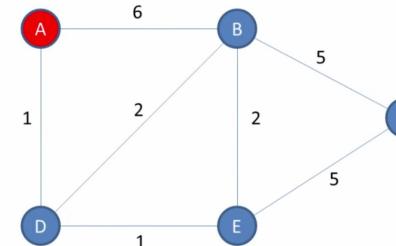
Until all vertices visited

```

10 # Single thread Dijkstra's implementation
11 def dijkstra(graph, initial):
12     """
13         GOAL: Shortest distances from node A to all other nodes
14         A node can be one of three states
15             Visited
16             Queue
17             Unvisited
18     """
19     visited = {initial: 0}
20     heap = [(0, initial)] # our priority queue
21     path = {}
22
23     nodes = set(graph.nodes) # unvisited state
24
25     while nodes and heap: #priority queue
26
27         current_weight, min_node = heapq.heappop(heap)
28         try:
29             while min_node not in nodes:
30                 current_weight, min_node = heapq.heappop(heap)
31         except IndexError:
32             break
33
34         nodes.remove(min_node)
35
36         for v in graph.edges[min_node]:
37             weight = current_weight + graph.distances[min_node, v]
38             if v not in visited or weight < visited[v]:
39                 visited[v] = weight
40                 heapq.heappush(heap, (weight, v))
41                 path[v] = min_node
42
43     return visited

```

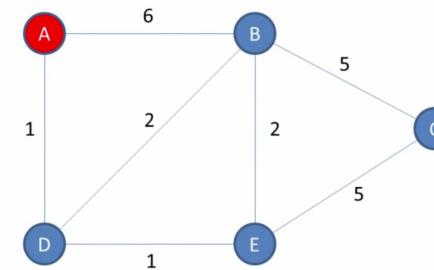
Consider the start vertex, A  
Distance to A from A = 0  
Distances to all other vertices from A are unknown, therefore  $\infty$  (infinity)



Vertex	Shortest distance from A	Previous vertex
A		
B		
C		
D		
E		

Visited = []      Unvisited = [A, B, C, D, E]

Visit the unvisited vertex with the smallest known distance from the start vertex  
*First time around, this is the start vertex itself, A*



Vertex	Shortest distance from A	Previous vertex
A	0	
B	$\infty$	
C	$\infty$	
D	$\infty$	
E	$\infty$	

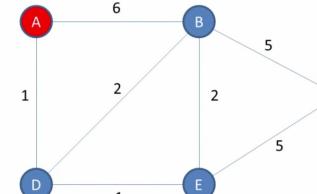
Visited = []      Unvisited = [A, B, C, D, E]

```

10 # Single thread Dijkstra's implementation
11 def dijkstra(graph, initial):
12     """
13         GOAL: Shortest distances from node A to all other nodes
14         A node can be one of three states
15             Visited
16             Queue
17             Unvisited
18     """
19     visited = {initial: 0}
20     heap = [(0, initial)] # our priority queue
21     path = {}
22
23     nodes = set(graph.nodes) # unvisited state
24
25     while nodes and heap: #priority queue
26
27         current_weight, min_node = heapq.heappop(heap)
28         try:
29             while min_node not in nodes:
29                 current_weight, min_node = heapq.heappop(heap)
30         except IndexError:
31             break
32
33         nodes.remove(min_node)
34
35         for v in graph.edges[min_node]:
36             weight = current_weight + graph.distances[min_node, v]
37             if v not in visited or weight < visited[v]:
38                 visited[v] = weight
39                 heapq.heappush(heap, (weight, v))
40                 path[v] = min_node
41
42     return visited

```

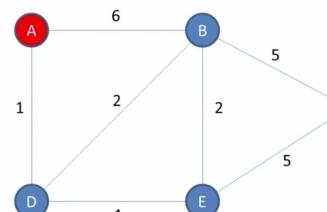
Consider the start vertex, A  
Distance to A from A = 0  
Distances to all other vertices from A are unknown, therefore  $\infty$  (infinity)



Visited = []      Unvisited = [A, B, C, D, E]

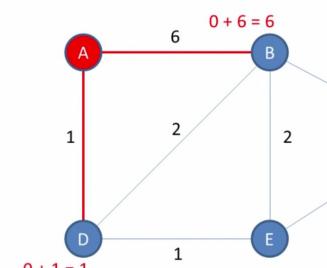
Visit the unvisited vertex with the smallest known distance from the start vertex  
First time around, this is the start vertex itself, A

Vertex	Shortest distance from A	Previous vertex
A		
B		
C		
D		
E		



For the current vertex, calculate the distance of each neighbour from the start vertex

Vertex	Shortest distance from A	Previous vertex
A	0	
B	$\infty$	
C	$\infty$	
D	$\infty$	
E	$\infty$	



Visited = []      Unvisited = [A, B, C, D, E]

Vertex	Shortest distance from A	Previous vertex
A	0	
B	$\infty$	
C	$\infty$	
D	$\infty$	
E	$\infty$	

# Algorithm

Let distance of start vertex from start vertex = 0

Let distance of all other vertices from start =  $\infty$  (infinity)

Repeat

- Visit the unvisited vertex with the smallest known distance from the start vertex

- For the current vertex, examine its unvisited neighbours

- For the current vertex, calculate distance of each neighbour from start vertex

- If the calculated distance of a vertex is less than the known distance, update the shortest distance

- Update the previous vertex for each of the updated distances

- Add the current vertex to the list of visited vertices

Until all vertices visited

# Algorithm

Greedy Algo: get to the end more quickly; makes locally optimal choices

Let distance of start vertex from start vertex = 0

Let distance of all other vertices from start =  $\infty$  (infinity)

Repeat

Visit the unvisited vertex with the smallest known distance from the start vertex

For the current vertex, examine its unvisited neighbours

For the current vertex, calculate distance of each neighbour from start vertex

If the calculated distance of a vertex is less than the known distance, update the shortest distance

Update the previous vertex for each of the updated distances

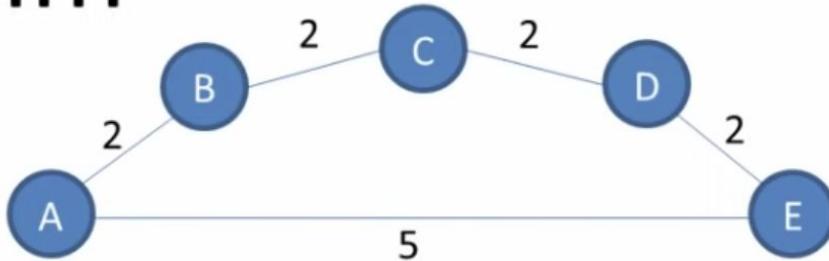
Add the current vertex to the list of visited vertices

Until all vertices visited

# Algorithm

Let distance of start vertex from start vertex = 0

Let distance of all other vertices from start =  $\infty$  (infinity)



Repeat

Visit the unvisited vertex with the smallest known distance from the start vertex

For the current vertex, examine its unvisited neighbours

For the current vertex, calculate distance of each neighbour from start vertex

If the calculated distance of a vertex is less than the known distance, update the shortest distance

Update the previous vertex for each of the updated distances

Add the current vertex to the list of visited vertices

Until all vertices visited

# Algorithm

Let distance of start vertex from start vertex = 0

Let distance of all other vertices from start =  $\infty$  (infinity)

WHILE vertices remain unvisited

    Visit unvisited vertex with smallest known distance from start vertex (call this 'current vertex')

    FOR each unvisited neighbour of the current vertex

        Calculate the distance from start vertex

        If the calculated distance of this vertex is less than the known distance

            Update shortest distance to this vertex

            Update the previous vertex with the current vertex

        end if

    NEXT unvisited neighbour

    Add the current vertex o the list of visited vertices

END WHILE

61-student-348823 &gt; w261

File Edit View Run Kernel Git Tabs Settings Help

DEMO9\_WORKBOOK\_MASTER.ipynb

1. MASTER Unit 9 - Graph Algorithms at Scale  
 1.0.1. Notebook Set-Up

2. Exercise 1. Graphs Overview

3. Warm-up questions:

4. Exercise 2. Data Structures Review.

5. Exercise 3. Graph Traversal.

5.1. Depth-First Search

5.2. Breadth-First Search

5.3. Other great algorithm books:

5.3.1. BFS

5.3.2. DFS

5.4. BFS versus DFS (no animation)

6. TASK: Write BFS Search Function In Apache Spark And Pandas Dataframe:

6.1. TODO: iterative search over undirected graph

7. Exercise 4. SSSP for unweighted graphs.

8. Exercise 5: Dijkstra's algorithm (SSSP for weighted graphs)

8.0.0.1. What about longest paths?

8.1. 8.1. ¶

9. Exercise 6. Distributed SSSP

9.0.0.1. Graph algorithms typically involve:

9.0.0.2. Key questions:

9.0.1. Distributed SSSP Algorithm

OLD\_demo8\_OLD\_workbookX demo8\_workbook\_MASTERF demo7\_workbook\_MASTERX hw3\_Workbook\_final\_mastX hw4\_Workbook\_MASTER.ipX

/ [1 files][ 0.0 B/ 0.0 B]  
 Operation completed over 1 objects.

[22]:

demo9\_workbook\_MASTERF c02w06\_BFS\_Spark\_SQL.ipX

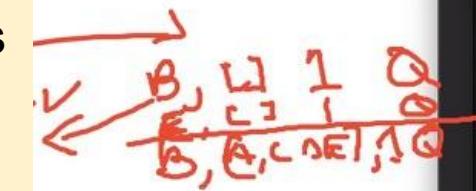
- Visited
- Queue
- Unvisited

9.0.1.2. Phase 1 - Initialize graph

- Start with Source node (node 1)

Dijkstra Alg.

See notebook for a sketch of Dijkstra's for weighted graphs using Spark Map-Reduce implementation



D	B,C,E	inf	U
E	A,B,D	inf	U

[12]:

```

1 # nodes
2 color_map = []
3 for node in G:
4     if node == "A":
5         color_map.append('#ff0055')
6     else:
7         color_map.append("#cccccc")
8 nx.draw_networkx_nodes(G, pos, node_size=1200, node_color=color_map)
9 nx.draw_networkx_edges(G, pos, edgelist=edges,
10                         width=1, alpha=0.5, edge_color='b')
11 # labels
12 nx.draw_networkx_labels(G, pos, font_size=12, font_family='sans-serif')
13 plt.axis('off')

```

Once a node is in V state, that record is finished ?? or isn't?

---

# Spark's GraphFrames library

## Databricks on AWS

[Get started](#)[What is Databricks?](#)[Tutorials and best practices](#)[Release notes](#)[Load data](#)[Explore data](#)[Prepare data](#)[Share data \(Delta sharing\)](#)[Data Marketplace](#)[Data engineering](#)

## Machine learning

[Try Databricks Machine Learning](#)[Prepare data and environment](#)[Train models](#)[Track model development](#)[Manage model lifecycle in Unity Catalog](#)[Manage model lifecycle using the Workspace Model Registry](#)[Serve and deploy models](#)[Export and import models](#)[Reference solutions](#)[Databricks Feature Store](#)[Documentation](#) > [Introduction to Databricks Machine Learning](#) > [GraphFrames](#) > GraphFrames user guide – Python

# GraphFrames user guide – Python

June 09, 2022

To learn more about GraphFrames, try importing and running this notebook in your workspace.

## GraphFrames Python notebook

[Open notebook in new tab](#) Copy link for import

## GraphFrames User Guide (Python)

This notebook demonstrates examples from the [GraphFrames User Guide](#).

The GraphFrames package is available from [Spark Packages](#).

```
from functools import reduce
from pyspark.sql.functions import col, lit, when
Expand notebook ▾ import *
```

[https://docs.databricks.com/\\_extras/notebooks/source/graphframes-user-guide-py.html](https://docs.databricks.com/_extras/notebooks/source/graphframes-user-guide-py.html)

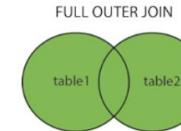
Was this article helpful?



**Tip:** FULL OUTER JOIN and FULL JOIN are the same.

## FULL OUTER JOIN Syntax

```
SELECT column_name(s)
FROM table1
FULL OUTER JOIN table2
ON table1.column_name = table2.column_name
WHERE condition;
```



**Note:** FULL OUTER JOIN can potentially return very large result-sets!

## Demo Database

In this tutorial we will use the well-known Northwind sample database.

Below is a selection from the "Customers" table:

CustomerID	CustomerName	ContactName	Address	City	PostalCode	Country
1	Alfreds Futterkiste	Maria Anders	Obere Str. 57	Berlin	12209	Germany
2	Ana Trujillo Emparedados y helados	Ana Trujillo	Avda. de la Constitución 2222	México D.F.	05021	Mexico
3	Antonio Moreno Taquería	Antonio Moreno	Mataderos 2312	México D.F.	05023	Mexico

And a selection from the "Orders" table:

OrderID	CustomerID	EmployeeID	OrderDate	ShipperID
10308	2	7	1996-09-18	3
10309	37	3	1996-09-19	1
10310	77	8	1996-09-20	2

## SQL FULL OUTER JOIN Example

The following SQL statement selects all customers, and all orders:

```
SELECT Customers.CustomerName, Orders.OrderID
FROM Customers
FULL OUTER JOIN Orders ON Customers.CustomerID=Orders.CustomerID
ORDER BY Customers.CustomerName;
```

A selection from the result set may look like this:

CustomerName	OrderID
Null	10309
Null	10310
Alfreds Futterkiste	Null
Ana Trujillo Emparedados y helados	10308
Antonio Moreno Taquería	Null

**Note:** The **FULL OUTER JOIN** keyword returns all matching records from both tables whether the other table matches or not. So, if there are rows in "Customers" that do not have matches in "Orders", or if there are rows in "Orders" that do not have matches in "Customers", those rows will be listed as well.

[https://www.w3schools.com/sql/sql\\_join\\_full.asp](https://www.w3schools.com/sql/sql_join_full.asp)